# AUTOMATED REASONING PROJECT

## REASONING ABOUT TRUTH

Graham Priest

# TECHNICAL REPORT SERIES



Research School of Social Sciences
Australian National University

TECHNICAL REPORT TR-ARP-2/88   22 February, 1988

# REASONING ABOUT TRUTH

Graham Priest

University of Western Australia

# Reasoning about Truth

## 0 Introduction

Any AI reasoning system with reasonable ambitions must have a way
of describing, specifying or representing situations, states of
affairs or wot not. Moreover, any AI reasoner that wants to
perform cognitive reasoning, a central aspect of people's
intelligence, must be able to express and reason about the
cognitive attitudes that are taken, by the reasoner and others,
to those representations: whether they are known, believed,
true, provable etc. (all quite distinct notions). Thus, for
example, A may reason 'What B has told me in the past has usually
been true. B now tells me that he has just heard from C. Hence
I may reasonably believe that B has heard from C.'

In the past several years, we have seen a number of studies
concerning the AI handling of cognitive reasoning. (See, for
example, any number of the papers in Halpern [86].) These
studies have concentrated on only some cognitive attitudes,
particularly, knowledge and belief. The notion of truth has been
largely ignored; and this despite the fact that it is, arguably,
one of the more central notions. For example, it is one of the
distinguishing characteristics between knowledge and belief; it
is a necessary condition of adequate proof; etc. The recent
paper by Perlis [85], which does address the problem of reasoning
about truth is therefore highly timely. Moreover, it creates yet
another bond between AI and logic; for reasoning about truth is
something that has formed a central part of logical
investigations this century.[1]

The following paper falls into two parts. In the first part I
will describe Perlis' construction and argue that it has certain
inadequacies. In the second part I will describe an approach to
the problem of reasoning about truth which I think is preferable,
and explain why. The approach draws on fairly recent work in a
branch of logic called 'paraconsistent logic'. This part of the
paper may therefore serve the function of introducing the reader
to parts of logic, relevant to AI, of which they may not be
aware.

## 1 Semantic Closure

The aim of an AI account of reasoning about truth (and related
notions) is to produce some formally tractable way of
representing the legitimate inferences that cognitive reasoners
are wont to make about truth (and its associated notions).
Obviously, the first prerequisite of such an account is to have a
language with a predicate 'is true' (which I will write as T).
Syntactically, the predicate is to apply to the representations
of states. We may conveniently take these to be sentences. This
is not only simple and apparently adequate, but is also the
dominant line that logicians have taken since Tarski, and so
allows the application of any established logical technology.

What, however, makes this predicate a *truth* predicate? If we
survey the inferences that characteristicaly involve the notion

of truth, we find essentially two. We infer '$\underline{\alpha}$ is true' from $\alpha$, and $\alpha$ from '$\underline{\alpha}$ is true' (where $\underline{\alpha}$ is a noun phrase which we can think of as a name for the the sentence $\alpha$). This suggests that the truth predicate, T, is characterised by the inference scheme:

$$T\underline{\alpha} \Leftrightarrow \alpha$$

which is called by logicians the T-scheme.[2]

It is perhaps rather surprising that such a trivial form of inference leads to trouble. Yet it does so. For all we need is a modicum of self-reference, obtainable in numerous ways, to find a state which claims that it, itself, is not true, i.e. a state, $\beta$, of the form $\neg T\underline{\beta}$. Applying the T-scheme to this, we we obtain:

$$T\underline{\beta} \Leftrightarrow \neg T\underline{\beta}$$

and hence:

$$T\underline{\beta} \wedge \neg T\underline{\beta}$$

This contradiction, the liar paradox, in itself, might not be too much of a problem. After all, is it surprising that such counter-intuitive results follow from consideration of such a pathological state? Unfortunately, if we then throw in the principle of standard logic that anything can be deduced from a contradiction a real problem arises. For the reasoner who has gone through this process can now infer everything. Which is slightly too much.

A dodge that logicians have used since the 30's to avoid this problem is to separate out the system to whose entities truth is attributed (the object language) and the system which attributes truth (the metalanguage), and claim that these must be distinct. (The construction is due to Tarski, though it should be said, in fairness to him, that he did not think that ordinary language reasoning about truth worked in this way.) Thus, the claim that this very state is not true can not be expressed at all. If it could be, then since it attributes truth, it must be in the metalanguage, but since it is that to which truth is attributed, it must be in the object language. Hence this is impossible.

In his paper Perlis argues, quite correctly, that this fix will not work. Cognitive representations are not intrinsically typed in this fashion, and any attempt to impose such a partition is not only artificial, but renders a great deal of perfectly correct and unproblematical reasoning impossible. This point is one of which logicians are now acutely aware, and most would agree that Perlis is quite right. As a result of this awareness, in recent years logicians have been investigating systems which aim at semantic closure, that is, systems that can talk about the truth of their own sentences. Characteristically, these approaches reject, or at least, weaken the T-scheme. Such approaches all face well known problems. (See for example Priest [84] and ch 1 of Priest [87], which also discuss the problems of the Tarskian approach further.) Perlis produces a novel such approach, which is not only simple, but works within the framework of orthodox logic. To this I now turn.

## 2 The T*-Scheme

Perlis' suggestion is simply to replace the T-scheme by what we will call the T*-scheme:[3]

$$T\underline{\alpha} \equiv \alpha^*  \qquad (T^*)$$

where $\alpha^*$ is the result of putting $\alpha$ into normal form (either conjunctive or disjunctive), and then replacing all occurrences of the form $\neg T\beta$ by $T(\neg\beta)$. Thus, in $\alpha^*$, only atomic sentences are negated, and atomic sentences containing the truth predicate are never negated. Since a formula is logically equivalent to its normal form, this preserves the T-scheme for those sentences that do not contain T, and even for those that do, provided that the occurrences of T are not within the scope of a negation.

The T*-scheme may look a little strange, and is certainly unlikely to occur to anyone *a priori*. The mystique may be removed, however, by considering its intuitive motivation.

As Perlis describes it, this builds on an idea of Kripke. First, we determine a class of sentences whose truth value may be fixed in a certain (transfinite) recursive fashion. These include (properly) all the sentences which do not contain T. Call these sentences *grounded*. The negation of a true grounded sentence is a false grounded sentence, and vice versa. On the Kripe construction, sentences that are not grounded are neither true nor false. The truth predicate applies truly to true sentences, falsely to false sentences, and neither-true-nor-falsely to ungrounded sentences. Consequently, $T\underline{\alpha}$ always has the same truth value (or lack of it) as $\alpha$.

As an analysis of truth, Kripke's construction is problematic for a number of reasons. One is that it does not dispense with the object language / metalanguage distinction. This is because the fact that a sentence is not true cannot be correctly expressed in the language itself. For if $\alpha$ is neither true nor false then $\neg T\underline{\alpha}$ is neither true nor false, not true as required. Moreover, not only does the T-scheme fail (if $\alpha$ is neither true nor false, so is the instance of the T-scheme for it), but no reasonable approximation to it seems to be available.

Perlis' suggestion is, in effect, to define a new classical interpretation of the language, $I$, such that the truth conditions of atomic sentences not containing T are the same as those in the Kripke interpretation, and those for sentences containing T are:

$T\underline{\alpha}$ is true in $I$ iff a is Kripke-true
$T\underline{\alpha}$ is false in $I$ otherwise.

It follows that all sentences are either true or false (in $I$). Moreover, all Kripke-true sentences are true (in $I$), and all Kripke-false sentences are false (as a simple induction shows). It therefore follows that:

$$T\underline{\alpha} \supset \alpha$$

is valid in $I$. The converse, however, is not. For if $\alpha$ is a true Kripke-neither sentence the consequent is true and the

antecedent is false.  Thus, the T-scheme, as is to be expected,
fails in general.  However, the T* scheme is valid.  I leave the
proof of this as an exercise for those familiar with Kripke's
construction.[4]  Thus, we see what semantics for the truth
predicate really underlie the T*-scheme.

To see what happens to the liar paradox in these semantics, note
that:

$$\neg(T\underline{\alpha} \wedge T\underline{\neg\alpha}) \qquad (1)$$

is true (in $I$) for every $\alpha$.  Next, note that if $\beta$ is the liar
sentence, the T*-scheme gives:

$$T\underline{\beta} \equiv T\underline{\neg\beta} \quad (=\beta^*)$$

Thus, by (1) and classical logic: $\neg T\underline{\beta}$, i.e., $\beta$.  The inference to
$T\underline{\beta}$ is, however, blocked.  Thus, the contradiction does not arise;
though we do have the rather odd $\beta \wedge \neg T\underline{\beta}$.  Moreover, $I$ provides a
consistent interpretation of the T*-scheme (and (1)), which
establishes that no other contradictions arise in the theory.[5]


## 3  Criticisms of this Account

Though Perlis' solution to the problem of how to formalise
reasoning concerning truth is neat, it will not work.  It is
wrong for both theoretical and practical reasons.  Let us start
with the theoretical reasons.

One objection to Perlis' construction is provided by the very
fact that the T-scheme does not hold in general.  There are a
number of arguments to the effect that the T-scheme must hold for
the truth predicate, that it, indeed, characterises truth.  Some
of the arguments are as ancient as Aristotle, and some as modern
as Frege.  I will not rehearse them here, since I do not wish
this to be a philosophical paper.  (Some of these arguments can
be found in Priest [87], sections 4.2, 4.3.)  Let us, therefore,
move on to more technical objections.

One of the weaknesses of Kripke's construction is that it does
not dispose of the object language / metalanguage distinction, as
I noted above.  But Perlis' construction is in exactly the same
boat.  For there is still no way in Perlis' construction of
expressing the fact that a sentence is true (in $I$)!  The easiest
way to see this is just to note that if the expressive power of
the language is sufficiently strong then, since the logic is
classical, we can apply Tarski's theorem to show that the set of
true sentences cannot be defined by any formula with one free
variable.  Thus Perlis' own talk of truth (in interpretation $I$)
must be conceived of as occurring within a distinct metalanguage;
and his claim to have gotten rid of a such a metalanguage (p 312)
is just false.

It follows, in particular, that the formula Tx, does not express
the claim that x is true.  In fact, as the semantics make clear,
$T\underline{\alpha}$ is true iff $\alpha$ is Kripke-true; but there are plenty of
sentences that are true but not Kripke-true.  As we noted, if $\beta$
is the liar sentence, $\beta \wedge \neg T\underline{\beta}$ is true (in $I$).  Thus, $\beta$ is one such

4

formula.⁶   Nor  is T even a good approximation to truth (in *I*),
since  some  of  the most fundamental facts about  truth  and  T
differ:   For  example,  for  every sentence either  it  or  its
negation  is true; but there are α's such that $T\underline{\alpha}vT\neg\underline{\alpha}$ is  false.
Similarly,  if α is not true then its negation is true; but there
are α's for which $\neg T\underline{\alpha} \supset T\neg\underline{\alpha}$ fails. (For counter-examples to both,
take α to be Kripke-neither.)

As we see,  Perlis' account is theoretically flawed.  It might be
suggested,  however,  that this doesn't matter since the point of
the  construction  is  not  a theoretical but  a  practical  one.
Specifically,  the  aim is to construct a formalisation that  can
represent our ordinary reasoning concerning truth; and, it may be
suggested, the T*-scheme is, in fact, adequate for this.  Indeed,
Perlis provides some nice examples of inferences involving the T-
scheme which are accounted for equally by the T*-scheme.

Unfortunately,  the theoretical inadequacies inevitably flow over
into practical ones.  Suppose, for example, that someone has the
job  of  having destroyed all and only those books  that  contain
some truth, i.e., they act on the command:

    $\exists x (Tx$ & book(y) & occurs_in(x y)$) \equiv$ destroy(y)

They  learn of book b that it contains inconsistent assertions on
pp 91 and 197.   They then reason that one of these must be true,
and hence that the book is to be destroyed.   (The  formalisation
of  this  is obvious.)  The situation might be  screwy,  but  the
reasoning  is  perfectly  sound and correct.   Yet it  cannot  be
represented in Perlis' approach,  just because,  as we noted  two
paragraphs back, $T\underline{\alpha}vT\neg\underline{\alpha}$ is not available.

Let me give another example,  which concerns the failure of the T-
scheme,  and which is a slight modification of one of Perlis' own.
Suppose  we  are given that anyone who speaks truly is  a  human.
(Vampires,  the  other  kind  of inhabitant of  Lower  Slobbovia,
always  lie.)  Two speakers,  Od and Id,  are heard to  speak  as
follows:

Id: Everything I say is not true.
Od: What Id says is not true.

We  can show that Od is human as follows.  Suppose that what  Id
says is true.   Then everything that Id says is not true.  Hence,
what Id says is not true.   Hence,  by *reductio*,  what Id says is
not true.   But Od said just that.   Hence he spoke truly.  He is
therefore  a  human.   I  leave a formalisation of  this  to  the
interested  reader.   The  important point to note is  just  that
having deduced that what Id says is not true (¬Tw), to then infer
that  Od spoke the truth (¬Tw ⊃ T¬Tw) is precisely an instance of
(the half of) the T-scheme that does not hold on Perlis' account.

For good measure,  we can also infer that Id is human.   We  have
established  that what Id said is not true.   Thus,  Id has  said
something true;  hence he is human.  Again, this reasoning cannot
be  represented  in  Perlis'  construction,   just  because  the
principle  $\neg T\underline{\alpha} \supset \neg\alpha$ fails.   Notice also,  that there is  nothing
problematic  about these reasonings due to  inconsistency.   The
situation is quite consistent.  (Suppose that Od and Id are human

and that Id has said (at least) one true thing.) There is
therefore no paradoxical "funny business".

We see that Perlis' construction does not allow for correct and
unproblematic cognitive reasoning about truth. Hence, it is not
only theoretically incorrect, but also practically inadequate.


## 4 A More Adequate Solution

I now wish to propose a more adequate solution. The T-scheme, we
have seen, must be part of any adequate representation of
cognitive reasoning. We have also seen that this gives rise to
contradictions. It would appear that this must be accepted.
What needs to be rejected is the view that everything may be
deduced from a contradiction. After all, the fact that
contradictions may arise in self-referential situations is not
particularly surprising, or even worrying. What *is* worrying (and
also surprising to someone who has not been indoctrinated by a
course on Frege/Russell logic) is that once a contradiction has
been inferred, everything follows: *ex contradictione quodlibet*.
If this rule fails then there is no reason why the contradictions
produced by the paradoxes of cognitive reasoning should not be
allowed to stand: they do no harm.[7]

Logics where *ex contradictione* fails are called *paraconsistent
logics*. These are now familiar to logicians, though to computer
scientists they may be less so. (Some uses have been made of
some of the more elementary relevant logics in the AI literature.
See, e.g. Levesque [84], Fagin and Halpern [85]. Relevant logics
are one kind of paraconsistent logic.) I shall not attempt a
review of such logics here. This can be found in Priest and
Routley [84] and Priest *et al* [88], ch 5.) Instead, I will
describe one of the simplest and most natural such logics, LP,
(see Priest [79], Priest [87], ch 5) and show how it can be
applied to the present situation.

LP is obtained by relaxing the classical assumption that
sentences cannot be both true and false. Thus, an interpretation
assigns to each atomic sentences one of the truth values {1}
(true and true only), {0} (false and false only), and {1,0}
(both). Truth conditions for non-atomic sentences are given in
the familiar classical way, except that truth and falsity, now
being independent, must each be considered. Thus, let us say
that $\alpha$ is true (under an interpretation) iff 1 is in its truth
value (under that interpretation); similarly, it is false iff 0
is in its truth value. Then, under an interpretation:

$\neg\alpha$ is true iff $\alpha$ is false
$\neg\alpha$ is false iff $\alpha$ is true

$\alpha\wedge\beta$ is true iff $\alpha$ is true and $\beta$ is true
$\alpha\wedge\beta$ is false iff $\alpha$ is false or $\beta$ is false.

Disjunction is treated dually. $\alpha\supset\beta$ is defined as $\neg\alpha\vee\beta$; $\alpha\equiv\beta$ is
defined as $(\alpha\supset\beta)\wedge(\beta\supset\alpha)$. Quantifiers, as in normal accounts, are
just thought of as (possibly infinitary) conjunctions and
disjunctions over the domain of interpretation. As can be
checked, these truth conditions are sufficient to give all

6

formulas one of the three truth values.  Logical consequence is defined in the standard way.  An interpretation is a *model* of a formula iff the formula is true in that interpretation; it is a model of a set of formulas iff it is a model of every formula in the set; and:

$$\Sigma \vDash \alpha \text{ iff every model of } \Sigma \text{ is a model of } \alpha$$

It is a simple job to show that this logic is paraconsistent. Take the evaluation that makes p both true and false, and q false only. This makes $p \wedge \neg p$ true (and false) and q not true. (LP might be more familiar to some people as Kleene's strong 3-valued logic with middle element designated.)

Let me mention, in passing, one variation on these semantics. This is obtained by allowing any subset of {1,0}, including the empty set, to be a truth value. Otherwise details are the same. These semantics are Dunn's semantics for Anderson and Belnap's system of zero degree entailment. (Discussed in Belnap [77], used by Levesque [84].) The main difference between these two systems is that the three valued system validates the law of excluded middle, $\alpha \vee \neg \alpha$, and indeed, all classical tautologies, whilst the four valued system has no logical truths. In the present context, I take this to be a distinct advantage for the three valued system. For the aim is to capture ordinary reasoning about truth; and the law of excluded middle is an integral part of much of this.


LP has a number of simple proof theories. (For example, a tableau system is given in Lin [86].) I will give a natural deduction system, sound and complete with respect to these semantics. This is obtained by modifying a standard natural deduction system for first order logic (that of Prawitz' [1965]). The modification is simply to replace the ordinary negation rules (¬I and ¬E) by:

$$\begin{array}{c} \overline{\alpha} \\ \cdot \\ \cdot \\ \cdot \\ \neg \beta \qquad \beta \\ \hline \neg \alpha \end{array} \text{CON} \qquad\qquad \begin{array}{c} \\ \hline \alpha \vee \neg \alpha \end{array} \text{LEM}$$

$$\begin{array}{c} \neg \neg \alpha \\ \hline \alpha \end{array} \text{DN}$$

where in CON, $\alpha$ is the only undischarged assumption, and no application of LEM occurs in its sub-proof. If we delete LEM we obtain a proof theory for zero degree entailment. If we drop the restriction on CON, we obtain classical logic.

Having got the background logic sorted out, to provide a system to reason about truth, we merely add the two rules:

$$\begin{array}{c} T\underline{\alpha} \\ \hline \alpha \end{array} \text{TE} \qquad\qquad \begin{array}{c} \alpha \\ \hline T\underline{\alpha} \end{array} \text{TI}$$

where α is a closed formula. Let us call this system of rules
TLP. As it stands, TLP is consistent (that is, no formula of the
form β∧¬β is provable.) This can be proved by noting that LP is
consistent, and then observing that TLP can be collapsed into LP
proofs merely be deleting T's and underlinings.) The consistency
is due, however, to the fact that, so far, no self-referential
machinery has been provided. As soon as this is provided,
inconsistency results. Thus, suppose we can produce a formula,
α, such that we can establish α≡¬T<u>α</u>; it is then a simple matter
to deduce α∧¬α. I leave this as an elementary exercise.

Although some contradictions are now provable, it would obviously
be disastrous if all were (i.e. if the system were trivial).
Fortunately, then, it can be shown that this is not the case. It
is possible to construct non-trivial TLP models of first order
arithmetic (which certainly contains enough self-referential
machinery), which show this. See Dowden [85]. In particular,
anything that is Kripke-false is not provable.


## 5 The Disjunctive Syllogism and Minimal Inconsistency

The inference engine TLP is not subject to the objections I
brought against Perlis' account. As may easily be checked, the
T-scheme: T<u>α</u>≡α is provable; and because of the T-rules the T
predicate defines the set of truths in any interpretation. Thus,
the account is not subject to Tarski's theorem concerning the
indefinability of truth.

There is, however, one important objection. Just because the
logic is paraconsistent, some inferences that are classically
valid are LP-invalid. *Ex contradictione quodlibet*, of course,
fails. However, this is well known to follow from simpler and
less intuitively puzzling inferences. One of these must
therefore have to fail. In fact, what fails is the disjunctive
syllogism:

    α ¬α∨β / β

detachment for material implication (sometimes called *modus
ponens*, though this name is appropriate only if ⊃ is a genuine
implication connective - something the very failure of detachment
gives grounds to doubt). In fact, the disjunctive syllogism is
the only classically valid inference to fail (in the sense that
if this is added to LP classical logic results). Yet it is
reasonable to object to my proposal that the failure shows its
inadequacy, since this inference is a part of our standard
reasoning - about truth or anything else. Indeed, both the
examples I gave in the previous section apply detachment to the
T-scheme.

There are two ways to meet this objection, both involving
extensions of the inferential machinery of LP. The simplest way
is as follows. Observe that to obtain an LP counter-example to
the disjunctive syllogism (or any other classically valid but LP-
invalid inference) we must render the situation inconsistent (by
making some formula both true and false). Now, it is both
plausible and natural to take consistency as a default

assumption. (For a defence of this see Priest [87], ch 9.) In that case it makes sense to implement a non-monotonic logic that implements this default, and which therefore allows the disjunctive syllogism provided that no pertinent inconsistency can be proved.

The simplest way of doing this is as follows. (I outline only the propositional case. Full details are given in Priest [88].) If $v$ is a propositional LP evaluation, let $v!$ be $\{p : p$ is a propositional parameter and $p \wedge \neg p$ is true under $v\}$. $v!$ is a measure of the inconsistency of an interpretation. Given a set of formulas, $\Sigma$, call $v$ a *minimally inconsistent* (*mi*) model of $\Sigma$ iff i) $v$ is a model of $\Sigma$, and ii) if $\mu!$ is properly contained in $v!$ then $\mu$ is not a model of $\Sigma$. That consistency is a default assumption means that we suppose there to be no more inconsistency than we are forced to suppose; and a natural way of making this idea precise is simply to restrict ourselves to mi models. Thus, define the default consequence relation $\vDash_m$ as follows:

$$\Sigma \vDash_m \alpha \text{ iff every mi model of } \Sigma \text{ is a model of } \alpha$$

This logic, $LP_m$, is non-monotonic and paraconsistent. (As may easily be checked $\{\neg p \vee q, p\} \vDash_m q$; but $\{\neg p \vee q, p, p \wedge \neg p\} \nvDash_m q$.) It extends LP, and gives all classical consequences if the premises are consistent. (See Priest [88] for proofs.) Hence, in consistent situations, the disjunctive syllogism and all other classical inferences are valid. In particular, both of the examples of section three (and all other examples where inconsistency does not rear its ugly head) can be represented in terms of $LP_m$, since these situations are consistent. Moreover, even in inconsistent situations, $LP_m$ still allows us to use the disjunctive syllogism provided only that the inconsistencies do not "get in the way". (Thus, for example, $\{p, \neg p \vee q, r \wedge \neg r\} \vDash_m q$.) Hence $LP_m$ validates all classical inferences except where inconsistency would make them naturally doubtful anyway.


## 6 Relevant Logic

The second way of meeting the objection is to extend the language of LP to include a genuine implication operator, $\rightarrow$, which satisfies (*inter alia*) detachment (*modus ponens*) — but not the principle $(\alpha \wedge \neg \alpha) \rightarrow \beta$. The T-scheme can now be formulated using this connective, and detachment from it becomes possible. The examples of section 3, for example, can be represented in this way.

A genuine implication operator can be added to LP in numerous ways. Relevant logicians, in particular, have studied how to give the semantics of such an operator; and LP can be embedded in relevance logics.[9] As I observed in section four, the semantics of LP are a fragment of the semantics of zero degree entailment. One possible approach is therefore to work with the extended semantics. This is unsatisfactory for two reasons. First, one looses all classical tautologies, such as the law of excluded middle (as I observed); second, and in any case, these semantics do not allow for nesting the connective $\rightarrow$, something one would surely want. It is better, therefore, to embed LP semantics in

those of a full relevant logic.

This is not the place to go into the semantics of relevant logics in detail. (Details can be found in Dunn [86] or Routley *et al* [82].) Let me, however, indicate one embedding. One kind of semantics for relevant logics is based on an algebraic structure of the form $\langle L, \wedge, \vee, *, \Rightarrow, F \rangle$, where $\langle L, \wedge, \vee, * \rangle$ is a De Morgan lattice, $F$ is a certain filter on the lattice and $\Rightarrow$ is a binary operation satisfing at least the condition: $a \leq b$ iff $a \Rightarrow b \in F$. An algebraic evaluation is a map from formulas into the lattice such that $\wedge$, $\vee$, $*$, and $\Rightarrow$ are the interpretations of $\wedge$, $\vee$, $\neg$ and $\rightarrow$, respectively. Semantic consequence is defined in terms of membership-of-$F$ preservation under all evaluations. Given any LP interpretation it is possible to construct such an algebra and embed the interpretation in it. Conversely, any such algebra can be cut down to an LP interpretation. This shows that LP is exactly the extensional (i.e., $\wedge$, $\vee$, $\neg$) fragment of the relevant logic. (Full details can be found in the appendix of Priest [80].)

It is worth observing that many theories based on relevant logics can be shown to be non-trivial even when the T-scheme is available. To see this, note that Brady [88] has shown a large class of relevant logics to be non-trivial (though inconsistent) when augmented by the abstraction scheme of naive set theory:

$$x \in \{y; \varphi\} \leftrightarrow \varphi(y/x) \qquad \text{Abs}$$

where / denotes substitution, and $x$ is free for $y$ in $\varphi$. Now, let $\underline{\alpha}$ be $\{x; \alpha\}$, where $x$ is the least variable, in some standard enumeration, not occurring in $\alpha$; and let $Tx$ be $\varphi \in x$. Then by Abs:

$$T\underline{\alpha} \leftrightarrow \varphi \in \{x; \alpha\} \leftrightarrow \alpha$$

Hence, any such logic can non-trivially model the T-scheme.


## 7 Final Observations

The last two sections explain different ways of extending LP so that suitable detachments are available. Which of these is preferable on a given occasion may depend on the context. Having a genuine implication connective will not take care of a detachment if the major premise cannot be expressed as a genuine conditional; $LP_m$ will (consistency permitting). But $LP_m$ will not allow one to express an indefeasible connection between $\alpha$ and $\beta$ (i.e., one where one can always get from $\alpha$ to $\beta$); having a genuine conditional will. Maybe, on occasions, it will be necessary to use both of these devices, though I have no example of this to offer. At any rate, I take it that, between them, they overcome the objection. Let me finish with three pertinent but miscellaneous comments.

a) The approach to reasoning about truth that I have advocated accepts the T-scheme and uses a paraconsistent logic to accommodate the consequent inconsistencies. It might be suggested that another possible line is to accept the T-scheme, but accommodate the inconsistencies via some other mechanism, for example, by applying truth-maintenance techniques. (See, e.g,

Doyle [79].) Thus, for example, starting with an instance of the T-scheme $T\underline{x}\equiv\neg T\underline{x}$ marked IN, the TMS would mark it OUT as soon as it noted that from it and it alone $\alpha\wedge\neg\alpha$ follows.

It would be hubris to claim that no approach like this can be made to work. However, any such approach based on classical logic faces a pretty devastating objection based on Curry paradoxes. (See Priest [80] or Priest [87], ch 6.) Suppose that the connective $\Rightarrow$ satisfies both detachment and absorption $(\alpha\Rightarrow(\alpha\Rightarrow\beta)$ / $\alpha\Rightarrow\beta)$. Suppose that we have suitable self-referential machinery, and thus, for an arbitrary formula, $\beta$, can construct a formula $T\underline{x}\Rightarrow\beta$ whose name is $\underline{x}$. The instance of the T-scheme for $\alpha$ is: $T\underline{x}\Leftrightarrow(T\underline{x}\Rightarrow\beta)$. Now, absorption gives $T\underline{x}\Rightarrow\beta$; whence detachment from right to left gives $T\underline{x}$. Putting these together, again by detachment, gives $\beta$. Thus an arbitrary formula follows from the T-scheme, even without the help of *ex contradictione*. Thus, running a TMS on a set of assumptions that includes the T-scheme would be just like running a TMS on the set of assumptions that contains all formulas. The results would be just as arbitrary, and just as meaningless. They would reflect nothing but the order of backtracking.

Just for the record, it is worth noting that LP and certain relevant logics do not fall foul of Curry paradoxes, as the non-triviality results cited above show. This is because LP does not validate detachment for material implication, and suitable relevant logics do not contain absorption (though some relevant logics do). There is as yet no non-triviality proof for $LP_m$ with the T-rules, but the Curry arguments certainly break down. Although $\{T\underline{x}\equiv(T\underline{x}\supset p)\}$ $\vdash_m p$, $\{T\underline{x}\equiv(T\underline{x}\supset p)$, $T\underline{\beta}\equiv(T\underline{\beta}\supset\neg p)\}$ gives neither $p$ nor $\neg p$ in $LP_m$.[9]

b) The second observation concerns other paradoxes of cognitive reasoning. It is not only truth that is known to lead to paradoxes, but plausible conditions on belief, knowledge, proof and other intensional operators similarly lead to contradictions. (See, e.g., Asher and Kamp [86], Thomason [86] for discussion and references.) It would take me too far afield in this paper to discuss these. But the fact that there is little agreement about how to handle them attests to the fact that all proposed solutions are problematic. Here I note only that these paradoxes in cognitive reasoning can be handled in exactly the same way as those concerning truth: we simply add the appropriate rules of proof for reasoning about knowledge, belief etc., and allow the contradictions to stand, since they can do no harm.

c) The final comment concerns the automated implementation of the systems I have described. Though it is simple enough to write algorithms for a number of these (for example, a proof-tree search will do for LP, and a model search will do for propositional $LP_m$) the problem of efficient algorithms remains to be investigated. Only for relevant logics has a start been made on this. Some Details of this can be found in Thistlewaite *et al* [86] and [87], and Bollen [86].[10]

# Notes

1. As will probably be clear from the following paper, I write from the logicians' side of this partition; however, I hope to succeed, at least partly, in crossing it.

2. Since various notions of implication will play a role in this paper, let me comment briefly on my notation. I will use $\Rightarrow$ as a generic implication connective (where its precise properties are not at issue); $\supset$ as material implication (always defined using negation and disjunction); and $\rightarrow$ as a *bona fide* implication guaranteed to satify at least *modus ponens*. Their respective bi-implications are $\Leftrightarrow$, $\equiv$ and $\leftrightarrow$.

3. In practice, Perlis also assumes other principles verified by the model construction to follow, such as (1) below, but which do not, as far as I can see, follow from T*.

4. Hint: First show the result for $\alpha$ in normal form. The only non-trivial part of this argument concerns negated T sentences; but here one can use the fact that $\neg T\beta$ is Kripke-true iff $T\neg\beta$ is Kripke-true. Next, observe that in the Kleene strong three valued logic $\alpha$ is equivalent to its normal form $\alpha'$. Hence, $T\underline{\alpha}$ is equivalent to $T\underline{\alpha'}$. Finally, observe that $\alpha'^*$ is just $\alpha^*$.

5. The model-construction in Perlis' paper, taken from Gilmore, is somewhat different, but the final model is the same. Again, the proof is not difficult to find, and I leave it as an exercise. Hint: show by induction that the extension of the truth predicate is the same at each level of the Kripke and Gilmore hierarchies.

6. Perlis, in effect, admits that his truth predicate just means Kripke-true: '...T[rue] is to be taken to mean Kripke's sense, i.e., grounded and true...' (p 312).

7. In fact, it has been argued quite independently of the paradoxes of cognitive reasoning, that inference engines suitable for reasoning from complex data should be paraconsistent. (See, e.g. Belnap [77].) For any but the most simplistic data bases and rule systems are liable to be inconsistent. Further, since there is no decision procedure for inconsistency, there is no general and effective way that the inconsistencies can be weeded out. We therefore have to live with them.

8. A suitable implication operator does not have to be relevant, however. See Priest [87], ch 6.

9. The following are mi-model counter-examples for p and $\neg$p respectively: p false only, $T\beta$ true only, $T\underline{\alpha}$ both; p true only, $T\underline{\alpha}$ true only, $T\beta$ both.

10. I would like to thank two anonymous referees of *Artificial Intelligence* for their helpful comments on a first draft of this paper. The paper was rewritten while I was a Project Visitor at the Automated Reasoning Project, Australian National University. I would also like to express my thanks to them.

## References

Asher N. and Kamp J. [86] 'The Knower's Paradox and Representational Theories of Attitudes' in Halpern [86].

Belnap N. [77] 'A Useful Four-Valued Logic' in *Modern Uses of Multiple-Valued Logic*, eds. J.M.Dunn and G.Epstein, Reidel, 1987.

Bollen A. [88] 'A Relevant Extension of PROLOG', research paper and implementation, Automated Reasoning Project, Australian National University, January 1988.

Brady R. [88] 'The Non-triviality of Dialectical Set Theory' in Priest *et al* [88].

Dowden B. [84] 'Accepting Inconsistencies from the Paradoxes', *Journal of Philosophical Logic* 13, 125-130, 1984.

Doyle J. [79] 'A Truth Maintenance System', *Artificial Intelligence* 12, 231-72, 1979.

Dunn J.M. [86] 'Relevance Logic and Entailment' in *Handbook of Philosophical Logic* Vol 3, ed. D. Gabbay and F. Gruender, Reidel, 1985.

Fagin R., and Halpern J. [85] 'Belief, Awareness and Limited Reasoning: preliminary report', *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, Morgan Kaufman, 1985.

Halpern J.Y. [86] *Theoretical Aspects of Reasoning about Knowledge*, Morgan Kaufman, 1985.

Levesque H. [84] 'A Logic of Implicit and Explicit Belief', *Proceedings of the National Conference on Artificial Intelligence*, Morgan Kaufman, 1984.

Lin F. [86] Tableau Systems for Some Paraconsistent Logics' to appear.

Perlis D. [85] 'Languages with Self-Reference 1: Foundations', *Artificial Intelligence* 25, 301-22, 1985.

Prawitz D. [65] *Proof Theory*, Almqvist and Wiksell, 1965.

Priest G. [79] 'Logic of Paradox', *Journal of Philosophical Logic* 8, 219-41, 1979.

Priest G. [80] 'Sense, Entailment and Modus Ponens', *Journal of Philosophical Logic* 9, 415-35, 1980.

Priest G. [84] 'Semantic Closure', *Studia Logica* 43, 117-129, 1984.

Priest G. [87] *In Contradiction*, Nijhoff.

Priest   G.   [88]   'Consistency   by   Default';   paper read   to   a
           seminar   at   the   Automated   Reasoning   Project,
           Australian National University,   January 1988;   to
           appear.

Priest   G.   and Routley   R.   [84]   'Introduction:   Paraconsistent
           Logics',   *Studia Logica* 43,   3-16,   1984.

Priest  G.,      Routley   R.   and  Norman J.     [88]   *Paraconsistent
           Logics,* Philosophia Verlag,   1988.

Routley R.   et   al,     [82]   *Relevant   Logics   and   their   Rivals,*
           Ridgeview,   1982.

Thistlewaite P.,   McRobbie   M.   and   Meyer R.   [86]   'The   KRIPKE
           Automated Theorem Proving System',   *Proc.   of   the
           Eighth Conference on Automated Reduction,*   ed.   J.
           Siekmann,   Springer,   1986.

Thistlewaite P.,   McRobbie   M.   and   Meyer   R.     [87]   *Automated
           Theorem-Proving in Non-Classical Logics',*   Wiley,
           1987.

Thomason   R.     [86]   'Paradoxes   and Semantic Representation'   in
           Halpern [86].

# Automated Reasoning Project
# Technical Report Series

Meyer, R.K. & Slaney, J; *Computing Demorgan Monoids*

Brink, C; *Some Background on Multisets*

Brink, C; *R¬-Algebras and R¬-model Structures as Power Constructs*

Brink, C. & Heidema, J; *A Verisimilar Ordering of Theories Phrased in a Propositional Language*

McRobbie, M.A., Meyer, R.K. & Thistlewaite, P.B.; *Towards Efficient Knowledge-Based Automated Theorem Proving for Non-Standard Logics*

Malkin, P.K. & Martin, E.P.; *Logical Matrix Generation and Testing*

Lavers, P.; *Relevance and Disjunctive Syllogism*

Martin, E.P.; *A Formalized Metalanguage for* P-W *and* S

Meyer, R.K. & Mortensen, C.; *Alien Intruders in Relevant Arithmetic*

Fuhrmann, A.; *Reflective Modalities and Theory Change*

Meyer, R.K., Giambrone, S., Urquhart, A., & Martin, E.P.; *Commutative Monoid and Semilattice Semantics for* P-W

McRobbie, M.A. & Siekmann, J.H.; *Artificial Intelligence: Perspesctives and Predictions*

Read, S.; *Proof Theory and Semantics for Relevant Logics*

Hazen, A.; *Predicative Arithmetic*

Bollen, A.W.; *A Relevant Extension to* PROLOG

Pelletier, J.; *Further Developments in* THINKER, *an Automated Theorem Prover: Proof Condensation, Adding Identity, 'Empirical' Issues in Computational Complexity, Pseudo-Parallel Subproof Development*

Siekmann, J.; *Unification Theory*

Siekmann, J. & Szabo, P.; *The Undecidability of the DA-Unification Problem*

Priest, G.; *Reasoning About Truth.*