# Circumscription—A Form of Non-Monotonic Reasoning

## John McCarthy

Stanford University, Stanford, CA, U.S.A.

#### ABSTRACT

Humans and intelligent computer programs must often jump to the conclusion that the objects they can determine to have certain properties or relations are the only objects that do. Circumscription formalizes such conjectural reasoning.

## 1. Introduction . The Qualification Problem

McCarthy [6] proposed a program with 'common sense' that would represent what it knows (mainly) by sentences in a suitable logical language. It would decide what to do by deducing a conclusion that it should perform a certain act. Performing the act would create a new situation, and it would again decide what to do. This requires representing both knowledge about the particular situation and general common sense knowledge as sentences of logic.

The 'qualification problem', immediately arose in representing general common sense knowledge. It seemed that in order to fully represent the conditions for the successful performance of an action, an impractical and implausible number of qualifications would have to be included in the sentences expressing them. For example, the successful use of a boat to cross a river requires, if the boat is a rowboat, that the oars and rowlocks be present and unbroken, and that they fit each other. Many other qualifications can be added, making the rules for using a rowboat almost impossible to apply, and yet anyone will still be able to think of additional requirements not yet stated.

Circumscription is a rule of conjecture that can be used by a person or program for 'jumping to certain conclusions'. Namely, the objects that can be shown to have a certain property P by reasoning from certain facts A are all the objects that satisfy P. More generally, circumscription can be used to conjecture that the tuples  $\langle x, y, ..., z \rangle$  that can be shown to satisfy a relation P(x, y, ..., z) are all the tuples satisfying this relation. Thus we circumscribe the set of relevant tuples. We can postulate that a boat can be used to cross a river unless 'something' prevents it. Then circumscription may be used to conjecture that the only entities that can prevent the use of the boat are those whose existence follows from the facts at hand. If no lack of oars or other circumstance preventing boat use is deducible, then the boat is concluded to be usable. The correctness of this conclusion depends on our having 'taken into account' all relevant facts when we made the circumscription.

Circumscription formalizes several processes of human informal reasoning. For example, common sense reasoning is ordinarily ready to jump to the conclusion that a tool can be used for its intended purpose unless something prevents its use. Considered purely extensionally, such a statement conveys no information; it seems merely to assert that a tool can be used for its intended purpose unless it can't. Heuristically, the statement is not just a tautologous disjunction; it suggests forming a plan to use the tool.

Even when a program does not reach its conclusions by manipulating sentences in a formal language, we can often profitably analyze its behavior by considering it to *believe* certain sentences when it is in certain states, and we can study how these *ascribed beliefs* change with time (see [9]). When we do such analyses, we again discover that successful people and programs must jump to such conclusions.

## 2. The Need for Non-Monotonic Reasoning

We cannot get circumscriptive reasoning capability by adding sentences to an axiomatization or by adding an ordinary rule of inference to mathematical logic. This is because the well known systems of mathematical logic have the following *monotonicity property*. If a sentence q follows from a collection A of sentences and  $A \subset B$ , then q follows from B. In the notation of proof theory: if  $A \vdash q$  and  $A \subset B$ , then  $B \vdash q$ . Indeed a proof from the premisses A is a sequence of sentences each of which is a either a premiss, an axiom or follows from a subset of the sentences occurring earlier in the proof by one of the rules of inference. Therefore, a proof from A can also serve as a proof from B. The semantic notion of entailment is also monotonic; we say that A entails q (written  $A \models q$ ) if q is true in all models of A. But if  $A \models q$  and  $A \subset B$ , then every model of B is also a model of A, which shows that  $B \models q$ .

Circumspection is a formalized rule of conjecture that can be used along with the rules of inference of first order logic. Predicate circumscription assumes that entities satisfy a given predicate only if they have to on the basis of a collection of facts. Domain circumscription conjectures that the 'known' entities are all there are. It turns out that domain circumscription, previously called minimal inference, can be subsumed under predicate circumscription.

We will argue using examples that humans use such 'non-monotonic' reasoning and that it is required for intelligent behavior. The default case reasoning of many computer programs [11] and the use of THNOT in MICROPLANNER [12] programs are also examples of non-monotonic reasoning, but possibly of a different kind from those discussed in this paper. Hewitt [5] gives the basic ideas of the PLANNER approach.

The result of applying circumscription to a collection A of facts is a sentence schema that asserts that the only tuples satisfying a predicate P(x, ..., z) are those whose doing so follows from the sentences of A. Since adding more sentences to Amight make P applicable to more tuples, circumscription is not monotonic. Conclusions derived from circumscription are conjectures that A includes all the relevant facts and that the objects whose existence follows from A are all the relevant objects.

A heuristic program might use circumscription in various ways. Suppose it circumscribes some facts and makes a plan on the basis of the conclusions reached. It might immediately carry out the plan, or be more cautious and look for additional facts that might require modifying it.

Before introducing the formalism, we informally discuss a well known problem whose solution seems to involve such non-monotonic reasoning.

## 3. Missionaries and Cannibals

The Missionaries and Cannibals puzzle, much used in AI, contains more than enough detail to illustrate many of the issues.

"Three missionaries and three cannibals come to a river. A rowboat that seats two is available. If the cannibals ever outnumber the missionaries on either bank of the river, the missionaries will be eaten. How shall they cross the river"

Obviously the puzzler is expected to devise a strategy of rowing the boat back and forth that gets them all across and avoids the disaster.

Amarel [1] considered several representations of the problem and discussed criteria whereby the following representation is preferred for purposes of AI, because it leads to the smallest state space that must be explored to find the solution. A state is a triple comprising the numbers of missionaries, cannibals and boats on the starting bank of the river. The initial state is 331, the desired final state is 000, and one solution is given by the sequence (331, 220, 321, 300, 311, 110, 221, 020, 031, 010, 021, 000).

We are not presently concerned with the heuristics of the problem but rather with the correctness of the reasoning that goes from the English statement of the problem to Amarel's state space representation. A generally intelligent computer program should be able to carry out this reasoning. Of course, there are the well known difficulties in making computers understand English, but suppose the English sentences describing the problem have already been rather directly translated into first order logic. The correctness of Amarel's representation is not an ordinary logical consequence of these sentences for two further reasons.

First, nothing has been stated about the properties of boats or even the fact that rowing across the river doesn't change the numbers of missionaries or cannibals or the capacity of the boat. Indeed it hasn't been stated that situations change as a result of action. These facts follow from common sense knowledge, so let us imagine that common sense knowledge, or at least the relevant part of it, is also expressed in first order logic.

The second reason we can't *deduce* the propriety of Amarel's representation is deeper. Imagine giving someone the problem, and after he puzzles for a while, he suggests going upstream half a mile and crossing on a bridge. "What bridge," you say. "No bridge is mentioned in the statement of the problem." And this dunce replies, "Well, they don't say there isn't a bridge." You look at the English and even at the translation of the English into first order logic, and you must admit that "they don't say" there is no bridge. So you modify the problem to exclude bridges and pose it again, and the dunce proposes a helicopter, and after you exclude that, he proposes a winged horse or that the others hang onto the outside of the boat while two row.

You now see that while a dunce, he is an inventive dunce. Despairing of getting him to accept the problem in the proper puzzler's spirit, you tell him the solution. To your further annoyance, he attacks your solution on the grounds that the boat might have a leak or lack oars. After you rectify that omission from the statement of the problem, he suggests that a sea monster may swim up the river and may swallow the boat. Again you are frustrated, and you look for a mode of reasoning that will settle his hash once and for all.

In spite of our irritation with the dunce, it would be cheating to put into the statement of the problem that there is no other way to cross the river than using the boat and that nothing can go wrong with the boat. A human doesn't need such an ad hoc narrowing of the problem, and indeed the only watertight way to do it might amount to specifying the Amarel representation in English. Rather we want to avoid the excessive qualification and get the Amarel representation by common sense reasoning as humans ordinarily do.

Circumscription is one candidate for accomplishing this. It will allow us to conjecture that no relevant objects exist in certain categories except those whose existence follows from the statement of the problem and common sense knowledge. When we *circumscribe* the first order logic statement of the problem together with the common sense facts about boats etc., we will be able to conclude that there is no bridge or helicopter. "Aha," you say, "but there won't be any oars either." No, we get out of that as follows: It is a part of common knowledge that a boat can be used to cross a river *unless there is something wrong with it or something else prevents using it*, and if our facts do not require that there be something that prevents crossing the river, circumscription will generate the conjecture that there isn't. The price is introducing as entities in our language the 'somethings' that may prevent the use of the boat.

If the statement of the problem were extended to mention a bridge, then the circumscription of the problem statement would no longer permit showing the non-existence of a bridge, i.e., a conclusion that can be drawn from a smaller collec-

tion of facts can no longer be drawn from a larger. This non-monotonic character of circumscription is just what we want for this kind of problem. The statement, "There is a bridge a mile upstream, and the boat has a leak." doesn't contradict the text of the problem, but its addition invalidates the Amarel representation.

In the usual sort of puzzle, there is a convention that there are no additional objects beyond those mentioned in the puzzle or whose existence is deducible from the puzzle and common sense knowledge. The convention can be explicated as applying circumscription to the puzzle statement and a certain part of common sense knowledge. However, if one really were sitting by a river bank and these six people came by and posed their problem, one wouldn't take the circumscription for granted, but one *would* consider the result of circumscription as a hypothesis. In puzzles, circumscription seems to be a rule of inference, while in life it is a rule of conjecture.

Some have suggested that the difficulties might be avoided by introducing probabilities. They suggest that the existence of a bridge is improbable. The whole situation involving cannibals with the postulated properties cannot be regarded as having a probability, so it is hard to take seriously the conditional probability of a bridge given the hypotheses. More to the point, we mentally propose to ourselves the normal non-bridge non-sea-monster interpretation *before* considering these extraneous possibilities, let alone their probabilities, i.e. we usually don't even introduce the sample space in which these possibilities are assigned whatever probabilities one might consider them to have. Therefore, regardless of our knowledge of probabilities, we need a way of formulating the normal situation from the statement of the facts, and non-monotonic reasoning seems to be required. The same considerations seem to apply to fuzzy logic.

Using circumscription requires that common sense knowledge be expressed in a form that says a boat can be used to cross rivers unless there is something that prevents its use. In particular, it looks like we must introduce into our ontology (the things that exist) a category that includes something wrong with a boat or a category that includes something that may prevent its use. Incidentally, once we have decided to admit something wrong with the boat, we are inclined to admit a lack of oars as such a something and to ask questions like, "Is a lack of oars all that is wrong with the boat?"

Some philosophers and scientists may be reluctant to introduce such *things*, but since ordinary language allows "*something wrong with the boat*" we shouldn't be hasty in excluding it. Making a suitable formalism is likely to be technically difficult as well as philosophically problematical, but we must try.

We challenge anyone who thinks he can avoid such entities to express in his favorite formalism, "Besides leakiness, there is something else wrong with the boat." A good solution would avoid counterfactuals as this one does.

Circumscription may help understand natural language, because if the use of natural language involves something like circumscription, it is understandable that the expression of general common sense facts in natural language will be difficult without some form of non-monotonic reasoning.

## 4. The Formalism of Circumscription

Let A be a sentence of first order logic containing a predicate symbol  $P(x_1, ..., x_n)$  which we will write  $P(\bar{x})$ . We write  $A(\Phi)$  for the result of replacing all occurrences of P in A by the predicate expression  $\Phi$ . (As well as predicate symbols, suitable  $\lambda$ -expressions are allowed as predicate expressions).

Definition The circumscription of P in A(P) is the sentence schema

$$4(\Phi) \land \forall \bar{x}. (\Phi(\bar{x}) \supset P(\bar{x})) \supset \forall \bar{x}. (P(\bar{x}) \supset \Phi(\bar{x})).$$
(1)

(1) can be regarded as asserting that the only tuples  $(\bar{x})$  that satisfy P are those that have to—assuming the sentence A. Namely, (1) contains a predicate parameter  $\Phi$  for which we may substitute an arbitrary predicate expression. (If we were using second order logic, there would be a quantifier  $\forall \Phi$  in front of (1).) Since (1) is an implication, we can assume both conjuncts on the left, and (1) lets us conclude the sentence on the right. The first conjunct  $A(\Phi)$  expresses the assumption that  $\Phi$ satisfies the conditions satisfied by P, and the second  $\forall \bar{x} . (\Phi(\bar{x}) \supset P(\bar{x}))$  expresses the assumption that the entities satisfying  $\Phi$  are a subset of those that satisfy P. The conclusion asserts the converse of the second conjunct which tells us that in this case,  $\Phi$  and P must coincide.

We write  $A \vdash_P q$  if the sentence q can be obtained by deduction from the result of circumscribing P in A. As we shall see  $\vdash_P$  is a non-monotonic form of inference, which we shall call *circumscriptive inference*.

A slight generalization allows circumscribing several predicates jointly; thus jointly circumscribing P and Q in A(P, Q) leads to

$$A(\Phi, \Psi) \land \forall \bar{x}. (\Phi(\bar{x}) \supset P(\bar{x})) \land \forall \bar{y}. (\Psi(\bar{y}) \supset Q(\bar{y}))$$
$$\supset \forall \bar{x}. (P(\bar{x}) \supset \Phi(\bar{x})) \land \forall \bar{y}. (O(\bar{y}) \supset \Psi(\bar{y}))$$
(2)

in which we can simultaneously substitute for  $\Phi$  and  $\Psi$ . The relation  $A \vdash_{P,Q} q$  is defined in a corresponding way. Although we do not give examples of joint circumscription in this paper, we believe it will be important in some AI applications.

Consider the following examples:

Example 1. In the blocks world, the sentence A may be

isblock 
$$A \wedge isblock B \wedge isblock C$$
 (3)

asserting that A, B and C are blocks. Circumscribing isblock in (3) gives the schema

$$\Phi(A) \land \Phi(B) \land \Phi(C) \land \forall x. (\Phi(x) \supset isblock \ x) \supset \forall x. (isblock \ x \supset \Phi(x)).$$
(4)

If we now substitute

$$\Phi(x) \equiv (x = A \lor x = B \lor x = C)$$
<sup>(5)</sup>

into (4) and use (3), the left side of the implication is seen to be true, and this gives

$$\forall x. (isblock \ x \supset (x = A \lor x = B \lor x = C)), \tag{6}$$

which asserts that the only blocks are A, B and C, i.e. just those objects that (3) requires to be blocks. This example is rather trivial, because (3) provides no way of generating new blocks from old ones. However, it shows that circumscriptive inference is non-monotonic since if we adjoin *isblock* D to (3), we will no longer be able to infer (6).

Example 2. Circumscribing the disjunction

$$isblock A \lor isblock B \tag{7}$$

leads to

$$(\Phi(A) \lor \Phi(B)) \land \forall x. (\Phi(x) \supset isblock \ x) \supset \forall x. (isblock \ x \supset \Phi(x)).$$
(8)

We may then substitute successively  $\Phi(x) \equiv (x = A)$  and  $\Phi(x) \equiv (x = B)$ , and these give respectively

$$(A = A \lor A = B) \land \forall x.(x = A \supset isblock \ x) \supset \forall x.(isblock \ x \supset x = A),$$
(9)

which simplifies to

$$isblock A \supset \forall x. (isblock x \supset x = A)$$
(10)

and

$$(B = A \lor B = B) \land \forall x. (x = B \supset isblock x) \supset \forall x. (isblock x \supset x = B),$$
(11)

which simplifies to

$$isblock \ B \supset \forall x. (isblock \ x \supset x = B).$$
(12)

(10), (12) and (7) yield

 $\forall x. (isblock \ x \supset x = A) \lor \forall x. (isblock \ x \supset x = B), \tag{13}$ 

which asserts that either A is the only block or B is the only block.

**Example 3.** Consider the following algebraic axioms for natural numbers, i.e., non-negative integers, appropriate when we aren't supposing that natural numbers are the only objects.

isnatnum 
$$0 \land \forall x.$$
 (isnatnum  $x \supset$  isnatnum succ x). (14)

Circumscribing isnatnum in (14) yields

$$\Phi(0) \land \forall x. (\Phi(x) \supset \Phi(succ x)) \land \forall x. (\Phi(x) \supset isnatnum x)$$
$$\supset \forall x. (isnatnum x \supset \Phi(x)).$$
(15)

(15) asserts that the only natural numbers are those objects that (14) forces to be natural numbers, and this is essentially the usual axiom schema of induction. We

can get closer to the usual schema by substituting  $\Phi(x) \equiv \Psi(x) \wedge isnatnum x$ . This and (14) make the second conjunct drop out giving

$$\Psi(0) \land \forall x.(\Psi(x) \supset \Psi(succ \ x)) \supset \forall x.(isnatnum \ x \supset \Psi(x)).$$
(16)

**Example 4.** Returning to the blocks world, suppose we have a predicate on(x, y, s) asserting that block x is on block y in situation s. Suppose we have another predicate above(x, y, s) which asserts that block x is above block y in situation s. We may write

$$\forall x \ y \ s.(on(x, y, s) \supset above(x, y, s))$$
(17)

and

$$\forall x \ y \ z \ s.(above(x, y, s) \land above(y, z, s) \supset above(x, z, s)), \tag{18}$$

i.e., above is a transitive relation. Circumscribing above in (17) and (18) gives

$$\forall x \ y \ s.(on(x, y, s) \supset \Phi(x, y, s)) \land \forall x \ y \ z \ s.(\Phi(x, y, s) \land \Phi(y, z, s) \supset \Phi(x, z, s)) \land \forall x \ y \ s.(\Phi(x, y, s) \supset above(x, y, s)) \supset \forall x \ y \ s.(above(x, y, s) \supset \Phi(x, y, s))$$
(19)

which tells us that above is the transitive closure of on.

In the preceding two examples, the schemas produced by circumscription play the role of axiom schemas rather than being just conjectures.

## 5. Domain Circumscription

The form of circumscription described in this paper generalizes an earlier version called *minimal inference*. Minimal inference has a semantic counterpart called *minimal entailment*, and both are discussed in [8] and more extensively in [3]. The general idea of minimal entailment is that a sentence q is minimally entailed by an axiom A, written  $A \models_m q$ , if q is true in all *minimal models* of A, where one model if is considered less than another if they agree on common elements, but the domain of the larger many contain elements not in the domain of the smaller. We shall call the earlier form *domain circumscription* to contrast it with the *predicate circumscription* discussed in this paper.

The domain circumscription of the sentence A is the sentence

$$Axiom(\Phi) \land A^{\Phi} \supset \forall x. \Phi(x).$$
<sup>(20)</sup>

where  $A^{\Phi}$  is the relativization of A with respect to  $\Phi$  and is formed by replacing each universal quantifier  $\forall x$ . in A by  $\forall x. \Phi(x) \supset$  and each existential quantifier  $\exists x$ . by  $\exists x. \Phi(x) \land Axiom(\Phi)$  is the conjunction of sentences  $\Phi(a)$  for each constant a and sentences  $\forall x. (\Phi(x) \supset \Phi(f(x)))$  for each function symbol f and the corresponding sentences for functions of higher arities. Domain circumscription can be reduced to predicate circumscription by relativizing A with respect to a new one place predicate called (say) *all*, then circumscribing *all* in  $A^{\text{all}} \wedge Axiom(all)$ , thus getting

$$Axiom(\Phi) \land A^{\Phi} \land \forall x.(\Phi(x) \supset all(x)) \supset \forall x.(all(x) \supset \Phi(x)).$$
(21)

Now we justify our using the name all by adding the axiom  $\forall x.all(x)$  so that (21) then simplifies precisely to (20).

In the case of the natural numbers, the domain circumscription of true, the identically true sentence, again leads to the axiom schema of induction. Here Axiom does all the work, because it asserts that 0 is in the domain and that the domain is closed under the successor operation.

## 6. The Model Theory of Predicate Circumscription

This treatment is similar to Davis's [3] treatment of domain circumscription. Pat Hayes [4] pointed out that the same ideas would work.

The intuitive idea of circumscription is saying that a tuple  $\bar{x}$  satisfies the predicate P only if it has to. It has to satisfy P if this follows from the sentence A. The model-theoretic counterpart of circumscription is *minimal entailment*. A sentence q is minimally entailed by A, iff q is true in all minimal models of A, where a model is minimal if as few as possible tuples  $\bar{x}$  satisfy the predicate P. More formally, this works out as follows.

**Definition.** Let M(A) and N(A) be models of the sentence A. We say that M is a submodel of N in P, writing  $M \leq_P N$ , if M and N have the same domain, all other predicate symbols in A besides P have the same extensions in M and N, but the extension of P in M is included in its extension in N.

**Definition.** A model M of A is called *minimal in P* iff  $M' \leq_P M$  only if M' = M. As discussed by Davis [3], minimal models do not always exist.

Definition. We say that A minimally entails q with respect to P, written  $A \models_P q$  provided q is true in all models of A that are minimal in P.

**Theorem.** Any instance of the circumscription of P in A is true in all models of A minimal in P, i.e., is minimally entailed by A in P.

**Proof.** Let M be a model of A minimal in P. Let P' be a predicate satisfying the left side of (1) when substituted for  $\Phi$ . By the second conjunct of the left side, P is an extension of P'. If the right side of (1) were not satisfied, P would be a proper extension of P'. In that case, we could get a proper submodel M' of M by letting M' agree with M on all predicates except P and agree with P' on P. This would contradict the assumed minimality of M.

Corollary. If  $A \vdash_P q$ , then  $A \models_P q$ .

While we have discussed minimal entailment in a single predicate P, the relation  $<_{P,Q}$ , models minimal in P and Q, and  $\models_{P,Q}$  have corresponding properties and a corresponding relation to the syntactic notion  $\vdash_{P,Q}$  mentioned earlier.

## 7. More on Blocks

The axiom

$$\forall x \ y \ s \ (\forall z \ \neg prevents(z, move(x, y), s) \supset on(x, y, result(move(x, y), s)))$$
(22)

states that unless something prevents it, x is on y in the situation that results from the action move(x, y).

We now list various 'things' that may prevent this action.

$$\forall x \ y \ s.(\neg isblock \ x \lor \neg isblock \ y \Rightarrow prevents(NONBLOCK, move(x, y), s))$$
(23)  
$$\forall x \ y \ s.(\neg clear(x, s) \lor \neg clear(y, s) \Rightarrow prevents(COVERED, move(x, y), s))$$
(24)  
$$\forall x \ y \ s.(tooheavy \ x \Rightarrow prevents(weight \ x, move(x, y), s)).$$
(25)

Let us now suppose that a heuristic program would like to move block A onto block C in a situation s0. The program should conjecture from (22) that the action move(A, C) would have the desired effect, so it must try to establish  $\forall z. \neg prevents(z, move(A, C), s0)$ . The predicate  $\lambda z. prevents(z, move(A, C), s0)$  can be circumscribed in the conjunction of the sentences resulting from specializing (23), (24) and (25), and this gives

$$(\neg isblock \ A \lor \neg isblock \ C \supset \Phi(NONBLOCK))$$
  

$$\land (\neg clear(A, s0) \lor \neg clear(C, s0) \supset \Phi(COVERED))$$
  

$$\land (tooheavy \ A \supset \Phi(weight \ A))$$
  

$$\land \forall z.(\Phi(z) \supset prevents(z, move(A, C), s0))$$
  

$$\supset \forall z.(prevents(z, move(A, C), s0) \supset \Phi(z))$$
(26)

which says that the only things that can prevent the move are the phenomena described in (23)–(25). Whether (26) is true depends on how good the program was in finding all the relevant statements. Since the program wants to show that nothing prevents the move, it must set  $\forall z. (\Phi(z) \equiv false)$ , after which (26) simplifies to

(isblock 
$$A \land$$
 isblock  $B \land$  clear( $A$ , s0)  $\land$  clear( $B$ , s0)  $\land \neg$  tooheavy  $A$   
 $\supset \forall z . \neg prevents(z, move(A, C), s0).$  (27)

We suppose that the premisses of this implication are to be obtained as follows:

### (1) isblock A and isblock B are explicitly asserted.

(2) Suppose that the only onness assertion explicitly given for situation s0 is on(A, B, s0). Circumscription of  $\lambda x y. on(x, y, s0)$  in this assertion gives

$$\Phi(A, B) \land \forall x \ y. (\Phi(x, y) \supset on(x, y, s0)) \supset \forall x \ y. (on(x, y, s0) \supset \Phi(x, y)),$$
(28)

and taking  $\Phi(x, y) \equiv x = A \land y = B$  yields

$$\forall x \ y.(on(x, y, s0) \supset x = A \land y = B).$$
<sup>(29)</sup>

Using

 $\forall x \ s.(clear(x, s) \equiv \forall y. \neg on(y, x, s))$ (30)

as the definition of *clear* yields the second two desired premisses.

(3)  $\neg$  tooheavy(x) might be explicitly present or it might also be conjectured by a circumscription assuming that if x were too heavy, the facts would establish it.

Circumscription may also be convenient for asserting that when a block is moved, everything that cannot be proved to move stays where it was. In the simple blocks world, the effect of this can easily be achieved by an axiom that states that all blocks except the one that is moved stay put. However, if there are various sentences that say (for example) that one block is attached to another, circumscription may express the heuristic situation better than an axiom.

#### 8. Remarks and Acknowledgments

(1) Circumscription is not a 'non-monotonic logic.' It is a form of non-monotonic reasoning augmenting ordinary first order logic. Of course, sentence schemata are not properly handled by most present general purpose resolution theorem provers. Even fixed schemata of mathematical induction when used for proving programs correct usually require human intervention or special heuristics, while here the program would have to use new schemata produced by circumscription. In [10] we treat some modalities in first order logic instead of in modal logic. In our opinion, it is better to avoid modifying the logic if at all possible, because there are many temptations to modify the logic, and it would be very difficult to keep them compatible.

(2) The default case reasoning provided in many systems is less general than circumscription. Suppose, for example, that a block x is considered to be on a block y only if this is explicitly stated, i.e., the default is that x is not on y. Then for each individual block x, we may be able to conclude that it isn't on block A, but we will not be able to conclude, as circumscription would allow, that there are no blocks on A. That would require a separate default statement that a block is clear unless something is stated to be on it.

(3) The conjunct  $\forall \bar{x}.(\Phi(\bar{x}) \supset P(\bar{x}))$  in the premiss of (1) is the result of suggestions by Ashok Chandra [2] and Patrick Hayes [4] whom I thank for their help. Without it, circumscribing a disjunction, as in the second example in Section 4, would lead to a contradiction. (4) The most direct way of using circumscription in AI is in a heuristic reasoning program that represents much of what it believes by sentences of logic. The program would sometimes apply circumscription to certain predicates in sentences. In particular, when it wants to perform an action that might be prevented by something, it circumscribes the prevention predicate in a sentence A representing the information being taken into account.

Clearly the program will have to include domain dependent heuristics for deciding what circumscriptions to make and when to take them back.

(5) In circumscription it does no harm to take irrelevant facts into account. If these facts do not contain the predicate symbol being circumscribed, they will appear as conjuncts on the left side of the implication unchanged. Therefore, the original versions of these facts can be used in proving the left side.

(6) Circumscription can be used in other formalisms than first order logic. Suppose for example that a set a satisfies a formula A(a) of set theory. The circumscription of this formula can be taken to be

$$\forall x. (A(x) \land (x \subset a) \supset (a \subset x)). \tag{31}$$

If a occurs in A(a) only in expressions of the form  $z \in a$ , then its mathematical properties should be analogous to those of predicate circumscription. We have not explored what happens if formulas like  $a \in z$  occur.

(7) The results of circumscription depend on the set of predicates used to express the facts. For example, the same facts about the blocks world can be axiomatized using the relation on or the relation above considered in Section 4 or also in terms of the heights and horizontal positions of the blocks. Since the results of circumscription will differ according to which representation is chosen, we see that the choice of representation has epistemological consequences if circumscription is admitted as a rule of conjecture. Choosing the set of predicates in terms of which to axiomatize as set of facts, such as those about blocks, is like choosing a co-ordinate system in physics or geography. As discussed in [9], certain concepts are definable only relative to a theory. What theory admits the most useful kinds of circumscription may be an important criterion in the choice of predicates. It may also be possible to make some statements about a domain like the blocks world in a form that does not depend on the language used.

(8) This investigation was supported in part by ARPA Contract MDA-903-76-C-0206, ARPA Order No. 2494, in part by NSF Grant MCS 78-00524, in part by the IBM 1979 Distinguished Faculty Program at the T. J. Watson Research Center, and in part by the Center for Advanced Study in the Behavioral Sciences.

#### REFERENCES

- Amarel, S., On representation of problems of reasoning about actions, in D. Michie (Ed.) Machine Intelligence 3 (Edinburgh University Press, Edinburgh, 1971), pp. 131-171.
- 2. Chandra, A., personal conversation (August 1979).
- 3. Davis, M., Notes on the mathematics of non-monotonic reasoning. Artificial Intelligence, this issue.

4. Hayes, P., personal conversation (September 1979).

- 5. Hewitt, C., Description and theoretical analysis (using schemata) of PLANNER: a language for proving theorems and manipulating models in a robot, MIT AI Laboratory TR-258 (1972).
- 6. McCarthy, J., Programs with common sense, in: Proceedings of the Teddington Conference on the Mechanization of Thought Processes (H.M. Stationery Office, London, 1960).
- 7. McCarthy, J., and Hayes, P., Some philosophical problems from the standpoint of Artificial Intelligence, in: D. Michie (Ed.), *Machine Intelligence* 4 (American Elsevier, New York, NY, 1969).
- 8. McCarthy, J., Epistemological problems of artificial intelligence, in: Proceedings of the Fifth International Joint Conference on Artificial Intelligence (MIT, Cambridge, MA, 1977).
- 9. McCarthy, J., Ascribing mental qualities to machines, in: M. Ringle (Ed.), Philosophical Perspectives in Artificial Intelligence (Harvester Press, July 1979).
- 10. McCarthy, J., First order theories of individual concepts and propositions, in: D. Michie (Ed.), Machine Intelligence 9 (University of Edinburgh Press, Edinburgh, 1979).
- 11. Reiter, R., A logic for default reasoning, Artificial Intelligence, this issue.
- 12. Sussman, G. J., Winograd, T., and Charniak, E., Micro-Planner Reference Manual, Al Memo 203, MIT AI Lab. (1971).

#### REASONING ABOUT KNOWLEDGE AND ACTION

Robert C. Moore Artificial Intelligence Laboratory Stanford University Stanford, California 94305

#### Abstract

This paper discusses the problems of representing and reasoning with information about knowledge and action. The first section discusses the importance of having systems that understand the concept of knowledge, and how knowledge is related to action. Section 2 points out some of the special problems that are involved in reasoning about knowledge, and section 3 presents a logic of knowledge based on the idea of possible worlds. Section 4 integrates this with a logic of actions and gives an example of reasoning in the combined system. Section 5 makes some concluding comments.

#### 1. Introduction

One of the most important concepts an intelligent system needs to understand is the concept of knowledge. AI systems need to understand what knowledge they and the systems or people they interact with have, what knowledge is needed to achieve particular goals, and how that knowledge can be obtained. This paper develops a formalism that provides a framework for stating and solving problems like these. For example, suppose that there is a safe that John wants to open. The common sense inferences that we would like to make might include:

If John knows the combination, he can immediately open the safe.

If John does not know the combination, he cannot immediately open the safe.

If John knows where the combination is written, he can read the combination and then open the safe.

In thinking about this example, consider how intimately the concept of knowledge is tied up with action. Reasoning about knowledge alone is of limited value. We may want to conclude from the fact that John knows A and B that he must also know C and D, but the real importance of such information is usually that it tells us something about what John can do or is likely to do. A major goal of my research has been to work out some of the interactions of knowing and doing.

That this area has received little attention in AI is somewhat surprising. It is frequently stated that good interactive AI programs will require good models of the people they are communicating with. Surely, one of the most important aspects of a model of another person is a model of what he knows. The only serious work on these problems in AI which I am aware of is a brief disscussion in McCarthy and Hayes (1969), and some more recent unpublished writings of McCarthy. In philosophy there is a substantial literature on the logic of knowledge and belief. A good introduction to this Is Hintikka (1962) and papers by Quine, Kaplan, and Hintikka in Linsky (1971). Many of the ideas I will use come from these papers.

In representing facts about knowledge and actions, I will use first-order predicate calculus, a practice which is currently unfathionable. It seems to be widely believed that use of predicate calculus necessarily leads to inefficient reasoning and information retrieval programs. I believe that this is an overreaction to earlier attempts to build domain-independent theorem provers based on resolution. More recent research, including my own M.S. thesis (Moore, 1975), suggests that predicate calculus can be treated in a more natural manner than resolution and combined with domain-dependent control information for greater efficiency. Furthermore, the problems of reasoning about knowledge seem to require the full ability to handle quantifiers and logical connectives which only predicate calculus posseses.

Section 2 of this paper attempts to bring out some of the special problems involved in reasoning about knowledge. Section 3 presents a formalism which I believe solves these problems, and Section 4 integrates this with a formalism for actions. Section 5 makes some concluding comments.

#### 2. Problems in Reasoning about Knowledge

Reasoning about knowledge presents special difficulties. It turns out that we cannot treat "know" as just another relation. If we can represent "Block1 is on Block2" by On(Block1,Block2), we might be tempted to represent "John knows that P" simply by Know(John,P). This approach glosses over a number of problems. We might be suspicious from the first, since P is not the name of an object but is rather a sentence (or proposition). The semantics of predicate calculus forbid the arbitrary intermingling of sentences and terms for good reason. For one thing, the second argument position of Know is a referentially opaque context. Ordinarily in logic we can freely substitute an expression for one that is extensionally equivalent (i.e., one that has the same referent or truth value), without affecting the truth of the formula that contains the expression. This is called referential transparency. For example, if X + Y = 7 and X = 3, then 3 + Y = 7. This pattern of reasoning is not valid with Know. We cannot infer from Know(John, (X + Y = 7)) and X = 3 that Know(John, (3 + Y = 7)) is true, since John might not know the value of X.

One possible solution to this problem is to make the second argument of Know the name of a formula rather than the formula itself. This is essentially the same idea as Goedel numbering, although it is not necessary to use such an obscure encoding as the natural numbers. We won't specify exactly how the encoding is done, but simply use "P" to represent a term denoting the formula P. The representation of "John knows that P" now becomes Know(John,"P(A)") We are no longer in any danger of infering Know(John,"P(A)") from Know(John,"P(B)") and A = B, because A is not contained in "P(A)". Only the name of A, i.e. "A", is contained, and since "A" does not equal "B", there is no problem.

There is, however, a more serious problem, the fact that people can reason with their knowledge. We would expect a reasoning system to have built into it the ability to conclude B from A and  $A \Rightarrow B$ . But if we treat Know as just an ordinary predicate, we will have no reason to suppose that Know(John,"A") and Know(John," $A \Rightarrow B$ ") might suggest Know(John,"B"). This problem is emphasised by the fact that there is no formal connection between a formula and its name. The fact that we