

What Is a Knowledge Representation?

Randall Davis, Howard Shrobe, and Peter Szolovits

■ Although knowledge representation is one of the central and, in some ways, most familiar concepts in AI, the most fundamental question about it—What is it?—has rarely been answered directly. Numerous papers have lobbied for one or another variety of representation, other papers have argued for various properties a representation should have, and still others have focused on properties that are important to the notion of representation in general.

In this article, we go back to basics to address the question directly. We believe that the answer can best be understood in terms of five important and distinctly different roles that a representation plays, each of which places different and, at times, conflicting demands on the properties a representation should have. We argue that keeping in mind all five of these roles provides a usefully broad perspective that sheds light on some long-standing disputes and can invigorate both research and practice in the field.

What is a knowledge representation? We argue that the notion can best be understood in terms of five distinct roles that it plays, each crucial to the task at hand:

First, a knowledge representation is most fundamentally a *surrogate*, a substitute for the thing itself, that is used to enable an entity to determine consequences by thinking rather than acting, that is, by reasoning about the world rather than taking action in it.

Second, it is a set of ontological commitments, that is, an answer to the question, In what terms should I think about the world?

Third, it is a fragmentary theory of intelligent reasoning expressed in terms of three components: (1) the representation's fundamental conception of intelligent reasoning, (2) the set of inferences that the representa-

tion sanctions, and (3) the set of inferences that it recommends.

Fourth, it is a medium for pragmatically efficient computation, that is, the computational environment in which thinking is accomplished. One contribution to this pragmatic efficiency is supplied by the guidance that a representation provides for organizing information to facilitate making the recommended inferences.

Fifth, it is a medium of human expression, that is, a language in which we say things about the world.

Understanding the roles and acknowledging their diversity has several useful consequences. First, each role requires something slightly different from a representation; each accordingly leads to an interesting and different set of properties that we want a representation to have.

Second, we believe the roles provide a framework that is useful for characterizing a wide variety of representations. We suggest that the fundamental mind set of a representation can be captured by understanding how it views each of the roles and that doing so reveals essential similarities and differences.

Third, we believe that some previous disagreements about representation are usefully disentangled when all five roles are given appropriate consideration. We demonstrate the clarification by revisiting and dissecting the early arguments concerning frames and logic.

Finally, we believe that viewing representations in this way has consequences for both research and practice. For research, this view provides one direct answer to a question of fundamental significance in the field. It also suggests adopting a broad perspective on

*a
representation
...
functions as
a surrogate
inside the
reasoner...*

what's important about a representation, and it makes the case that one significant part of the representation endeavor—capturing and representing the richness of the natural world—is receiving insufficient attention. We believe that this view can also improve practice by reminding practitioners about the inspirations that are the important sources of power for a variety of representations.

Terminology and Perspective

Two points of terminology assist our presentation. First, we use the term *inference* in a generic sense to mean any way to get new expressions from old. We rarely talk about sound logical inference and, when doing so, refer to it explicitly.

Second, to give them a single collective name, we refer to the familiar set of basic representation tools, such as logic, rules, frames, and semantic nets, as knowledge representation technologies.

It also proves useful to take explicit note of the common practice of building knowledge representations in multiple levels of languages, typically, with one of the knowledge representation technologies at the bottom level. Hayes's (1978) ontology of liquids, for example, is at one level a representation composed of concepts like pieces of space, with portals, faces, sides, and so on. The language at the next, more primitive (and, as it turns out, bottom) level is first-order logic, where, for example, $In(s_1, s_2)$ is a relation expressing that space s_1 is contained in s_2 .

This view is useful in part because it allows our analysis and discussion to concentrate largely on the knowledge representation technologies. As the primitive representational level at the foundation of knowledge representation languages, those technologies encounter all the issues central to knowledge representation of any variety. They are also useful exemplars because they are widely familiar to the field, and there is a substantial body of experience with them to draw on.

What Is a Knowledge Representation?

Perhaps the most fundamental question about the concept of knowledge representation is, What is it? We believe that the answer is best understood in terms of the five fundamental roles that it plays.

Role 1: A Knowledge Representation Is a Surrogate

Any intelligent entity that wants to reason about its world encounters an important, inescapable fact: Reasoning is a process that goes on internally, but most things it wants to reason about exist only externally. A program (or person) engaged in planning the assembly of a bicycle, for example, might have to reason about entities such as wheels, chains, sprockets, and handle bars, but such things exist only in the external world.

This unavoidable dichotomy is a fundamental rationale and role for a representation: It functions as a surrogate inside the reasoner, a stand-in for the things that exist in the world. Operations on and with representations substitute for operations on the real thing, that is, substitute for direct interaction with the world. In this view, reasoning itself is, in part, a surrogate for action in the world when we cannot or do not (yet) want to take that action.¹

Viewing representations as surrogates leads naturally to two important questions. The first question about any surrogate is its intended identity: What is it a surrogate for? There must be some form of correspondence specified between the surrogate and its intended referent in the world; the correspondence is the semantics for the representation.

The second question is fidelity: How close is the surrogate to the real thing? What attributes of the original does it capture and make explicit, and which does it omit? Perfect fidelity is, in general, impossible, both in practice and in principle. It is impossible in principle because any thing other than the thing itself is necessarily different from the thing itself (in location if nothing else). Put the other way around, the only completely accurate representation of an object is the object itself. All other representations are inaccurate; they inevitably contain simplifying assumptions and, possibly, artifacts.

Two minor elaborations extend this view of representations as surrogates. First, it appears to serve equally well for intangible objects as well as tangible objects such as gear wheels: Representations function as surrogates for abstract notions such as actions, processes, beliefs, causality, and categories, allowing them to be described inside an entity so it can reason about them. Second, formal objects can of course exist inside the machine with perfect fidelity: Mathematical entities, for example, can be captured exactly, precisely because they are formal objects. Because almost any reasoning task will

encounter the need to deal with *natural objects* (that is, those encountered in the real world) as well as formal objects, imperfect surrogates are pragmatically inevitable.

Two important consequences follow from the inevitability of imperfect surrogates. One consequence is that in describing the natural world, we must inevitably lie, by omission at least. At a minimum, we must omit some of the effectively limitless complexity of the natural world; in addition, our descriptions can introduce artifacts not present in the world.

The second and more important consequence is that all sufficiently broad-based reasoning about the natural world must eventually reach conclusions that are incorrect, independent of the reasoning process used and independent of the representation employed. Sound reasoning cannot save us: If the world model is somehow wrong (and it must be), some conclusions will be incorrect, no matter how carefully drawn. A better representation cannot save us: All representations are imperfect, and any imperfection can be a source of error.

The significance of the error can, of course, vary; indeed, much of the art of selecting a good representation is in finding one that minimizes (or perhaps even eliminates) error for the specific task at hand. But the unavoidable imperfection of surrogates means that we can supply at least one guarantee for any entity reasoning in any fashion about the natural world: If it reasons long enough and broadly enough, it is guaranteed to err.

Thus, drawing only sound inferences does not free reasoning from error; it can only ensure that inference is not the source of the error. Given that broad-based reasoning is inevitably wrong, the step from sound inference to other models of inference is thus not a move from total accuracy to error, but is instead a question of balancing the possibility of one more source of error against the gains (for example, efficiency) it might offer.

We do not suggest that unsound reasoning ought to be embraced casually, but we do claim that given the inevitability of error, even with sound reasoning, it makes sense to pragmatically evaluate the relative costs and benefits that come from using both sound and unsound reasoning methods.

Role 2: A Knowledge Representation Is a Set of Ontological Commitments

If, as we argue, all representations are imperfect approximations to reality, each approximation attending to some things and ignoring others, then in selecting any repre-

sentation, we are in the very same act unavoidably making a set of decisions about how and what to see in the world. That is, selecting a representation means making a set of ontological commitments.² The commitments are, in effect, a strong pair of glasses that determine what we can see, bringing some part of the world into sharp focus at the expense of blurring other parts.

These commitments and their focusing-blurring effect are not an incidental side effect of a representation choice; they are of the essence: A knowledge representation is a set of ontological commitments. It is unavoidably so because of the inevitable imperfections of representations. It is usefully so because judicious selection of commitments provides the opportunity to focus attention on aspects of the world that we believe to be relevant.

The focusing effect is an essential part of what a representation offers because the complexity of the natural world is overwhelming. We (and our reasoning machines) need guidance in deciding what in the world to attend to and what to ignore. The glasses supplied by a representation can provide this guidance: In telling us what and how to see, they allow us to cope with what would otherwise be untenable complexity and detail. Hence, the ontological commitment made by a representation can be one of its most important contributions.

There is a long history of work attempting to build good ontologies for a variety of task domains, including early work on an ontology for liquids (Hayes 1978), the lumped element model widely used in representing electronic circuits (for example, Davis and Shrobe [1983]) as well as ontologies for time, belief, and even programming itself. Each of these ontologies offers a way to see some part of the world.

The lumped-element model, for example, suggests that we think of circuits in terms of components with connections between them, with signals flowing instantaneously along the connections. This view is useful, but it is not the only possible one. A different ontology arises if we need to attend to the electro-dynamics in the device: Here, signals propagate at finite speed, and an object (such as a resistor) that was previously viewed as a single component with an input-output behavior might now have to be thought of as an extended medium through which an electromagnetic wave flows.

Ontologies can, of course, be written down in a wide variety of languages and notations

All representations are imperfect, and any imperfection can be a source of error.

(for example, logic, Lisp); the essential information is not the form of this language but the *content*, that is, the set of concepts offered as a way of thinking about the world. Simply put, the important part is notions such as connections and components, and not whether we choose to write them as predicates or Lisp constructs.

The commitment we make by selecting one or another ontology can produce a sharply different view of the task at hand. Consider the difference that arises in selecting the lumped element view of a circuit rather than the electrodynamic view of the same device. As a second example, medical diagnosis viewed in terms of rules (for example, MYCIN) looks substantially different from the same task viewed in terms of frames (for example, INTERNIST). Where MYCIN sees the medical world as made up of empirical associations connecting symptom to disease, INTERNIST sees a set of prototypes, in particular prototypical diseases, that are to be matched against the case at hand.

Commitment Begins with the Earliest Choices The INTERNIST example also demonstrates that there is significant and unavoidable ontological commitment even at the level of the familiar representation technologies. Logic, rules, frames, and so on, embody a viewpoint on the kinds of things that are important in the world. Logic, for example, involves a (fairly minimal) commitment to viewing the world in terms of individual entities and relations between them. Rule-based systems view the world in terms of attribute-object-value triples and the rules of plausible inference that connect them, while frames have us thinking in terms of prototypical objects.

Thus, each of these representation technologies supplies its own view of what is important to attend to, and each suggests, conversely, that anything not easily seen in these terms may be ignored. This suggestion is, of course, not guaranteed to be correct because anything ignored can later prove to be relevant. But the task is hopeless in principle—every representation ignores something about the world; hence, the best we can do is start with a good guess. The existing representation technologies supply one set of guesses about what to attend to and what to ignore. Thus, selecting any of them involves a degree of ontological commitment: The selection will have a significant impact on our perception of, and approach to, the task and on our perception of the world being modeled.

The Commitments Accumulate in Layers

The ontological commitment of a representation thus begins at the level of the representation technologies and accumulates from there. Additional layers of commitment are made as we put the technology to work. The use of framelike structures in INTERNIST illustrates. At the most fundamental level, the decision to view diagnosis in terms of frames suggests thinking in terms of prototypes, defaults, and a taxonomic hierarchy. But what are the prototypes of, and how will the taxonomy be organized?

An early description of the system (Pople 1982) shows how these questions were answered in the task at hand, supplying the second layer of commitment:

The knowledge base underlying the INTERNIST system is composed of two basic types of elements: disease entities and manifestations.... [It] also contains a ... hierarchy of disease categories, organized primarily around the concept of organ systems, having at the top level such categories as "liver disease," "kidney disease," etc. (pp. 136–137)

Thus, the prototypes are intended to capture prototypical diseases (for example, a classic case of a disease), and they will be organized in a taxonomy indexed around organ systems. This set of choices is sensible and intuitive, but clearly, it is not the only way to apply frames to the task; hence, it is another layer of ontological commitment.

At the third (and, in this case, final) layer, this set of choices is instantiated: Which diseases will be included, and in which branches of the hierarchy will they appear? Ontological questions that arise even at this level can be fundamental. Consider, for example, determining which of the following are to be considered *diseases* (that is, abnormal states requiring cure): alcoholism, homosexuality, and chronic fatigue syndrome. The ontological commitment here is sufficiently obvious and sufficiently important that it is often a subject of debate in the field itself, independent of building automated reasoners.

Similar sorts of decisions have to be made with all the representation technologies because each of them supplies only a first-order guess about how to see the world: They offer a way of seeing but don't indicate how to instantiate this view. Frames suggest prototypes and taxonomies but do not tell us which things to select as prototypes, and rules suggest thinking in terms of plausible inferences but don't tell us which plausible inferences to attend to. Similarly, logic tells us to view the world in terms of individuals

and relations but does not specify which individuals and relations to use. Thus, commitment to a particular view of the world starts with the choice of a representation technology and accumulates as subsequent choices are made about how to see the world in these terms.

Reminder: A Knowledge Representation Is Not a Data Structure Note that at each layer, even the first (for example, selecting rules or frames), the choices being made are about representation, not data structures. Part of what makes a language representational is that it carries meaning (Hayes 1979; Brachman and Levesque 1985); that is, there is a correspondence between its constructs and things in the external world. In turn, this correspondence carries with it a constraint.

A semantic net, for example, is a representation, but a graph is a data structure. They are different kinds of entity, even though one is invariably used to implement the other, precisely because the net has (should have) a semantics. This semantics will be manifest in part because it constrains the network topology: A network purporting to describe family memberships as we know them cannot have a cycle in its parent links, but graphs (that is, data structures) are, of course, under no such constraint and can have arbitrary cycles.

Although every representation must be implemented in the machine by some data structure, the representational property is in the correspondence to something in the world and in the constraint that correspondence imposes.

Role 3: A Knowledge Representation Is a Fragmentary Theory of Intelligent Reasoning

The third role for a representation is as a fragmentary theory of intelligent reasoning. This role comes about because the initial conception of a representation is typically motivated by some insight indicating how people reason intelligently or by some belief about what it means to reason intelligently at all.

The theory is fragmentary in two distinct senses: (1) the representation typically incorporates only part of the insight or belief that motivated it and (2) this insight or belief is, in turn, only a part of the complex and multifaceted phenomenon of intelligent reasoning.

A representation's theory of intelligent reasoning is often implicit but can be made more evident by examining its three components: (1) the representation's fundamental conception of intelligent inference, (2) the set of inferences that the representation sanc-

tions, and (3) the set of inferences that it recommends.

Where the sanctioned inferences indicate what can be inferred at all, the recommended inferences are concerned with what should be inferred. (Guidance is needed because the set of sanctioned inferences is typically far too large to be used indiscriminately.) Where the ontology we examined earlier tells us how to see, the recommended inferences suggest how to reason.

These components can also be seen as the representation's answers to three corresponding fundamental questions: (1) What does it mean to reason intelligently? (2) What can we infer from what we know? and (3) What should we infer from what we know? Answers to these questions are at the heart of a representation's spirit and mind set; knowing its position on these issues tells us a great deal about it.

We begin with the first of these components, examining two of several fundamentally different conceptions of intelligent reasoning that have been explored in AI. These conceptions and their underlying assumptions demonstrate the broad range of views on the question and set important context for the remaining components.

What Is Intelligent Reasoning? What are the essential, defining properties of intelligent reasoning? As a consequence of the relative youth of AI as a discipline, insights about the nature of intelligent reasoning have often come from work in other fields. Five fields—mathematical logic, psychology, biology, statistics, and economics—have provided the inspiration for five distinguishable notions of what constitutes intelligent reasoning (table 1).

One view, historically derived from mathematical logic, makes the assumption that intelligent reasoning is some variety of formal calculation, typically deduction; the modern exemplars of this view in AI are the logicians. A second view, rooted in psychology, sees reasoning as a characteristic human behavior and has given rise to both the extensive work on human problem solving and the large collection of knowledge-based systems.

A third approach, loosely rooted in biology, takes the view that the key to reasoning is the architecture of the machinery that accomplishes it; hence, reasoning is a characteristic stimulus-response behavior that emerges from the parallel interconnection of a large collection of very simple processors. Researchers working on several varieties of connectionism are the current descendants of this line of

Mathematical Logic	Psychology	Biology	Statistics	Economics
Aristotle				
Descartes				
Boole	James		Laplace	Bentham Pareto
Frege			Bernoullii	Friedman
Peano	Hebb	Lashley	Bayes	
Goedel	Bruner	Rosenblatt		
Post	Miller	Ashby	Tversky,	Von Neumann
Church	Newell,	Lettvin	Kahneman	Simon
Turing	Simon	McCulloch, Pitts		Raiffa
Davis		Heubel, Weisel		
Putnam				
Robinson				
Logic PROLOG	SOAR KBS, Frames	Connectionism	Causal Networks	Rational Agents

Table 1. Views of Intelligent Reasoning and Their Intellectual Origins.

work. A fourth approach, derived from probability theory, adds to logic the notion of uncertainty, yielding a view in which *reasoning intelligently* means obeying the axioms of probability theory. A fifth view, from economics, adds the further ingredient of values and preferences, leading to a view of intelligent reasoning that is defined by adherence to the tenets of utility theory.

Briefly exploring the historical development of the first two of these views (the logical and the psychological) illustrates the different conceptions they have of the fundamental nature of intelligent reasoning and demonstrates the deep-seated differences in mind set that arise as a consequence.

Consider first the tradition that surrounds mathematical logic as a view of intelligent reasoning. This view has its historical origins in Aristotle's efforts to accumulate and catalogue the syllogisms in an attempt to determine what should be taken as a convincing argu-

ment.³ The line continues with René Descartes, whose analytic geometry showed that Euclid's work, apparently concerned with the stuff of pure thought (lines of zero width, perfect circles of the sorts only the gods could make), could, in fact, be married to algebra, a form of calculation, something mere mortals can do.

By the time of Gottfried Wilhelm von Leibnitz in the seventeenth century, the agenda was specific and telling: He sought nothing less than a *calculus of thought*, one that would permit the resolution of all human disagreement with the simple invocation, "Let us compute." By this time, there was a clear and concrete belief that as Euclid's once godlike and unreachable geometry could be captured with algebra, so some (or perhaps any) variety of that ephemeral stuff called thought might be captured in calculation, specifically, logical deduction.

In the nineteenth century, G. Boole provid-

ed the basis for propositional calculus in his “Laws of Thought”; later work by G. Frege and G. Peano provided additional foundation for the modern form of predicate calculus. Work by M. Davis, H. Putnam, and G. Robinson in the twentieth century provides the final steps in sufficiently mechanizing deduction to enable the first automated theorem provers. The modern offspring of this line of intellectual development include the many efforts that use first-order logic as a representation and some variety of deduction as the reasoning engine as well as the large body of work with the explicit agenda of making logical reasoning computational, exemplified by PROLOG.

This line of development clearly illustrates how approaches to representation are founded on and embed a view of the nature of intelligent reasoning. There is here, for example, the historical development of the underlying premise that reasoning intelligently means reasoning logically; anything else is a mistake or an aberration. Allied with this premise is the belief that *logically*, in turn, means first-order logic, typically, sound deduction. By simple transitivity, these two theories collapse into one key part of the view of intelligent reasoning underlying logic: Reasoning intelligently means reasoning in the fashion defined by first-order logic. A second important part of the view is the allied belief that intelligent reasoning is a process that can be captured in a formal description, particularly a formal description that is both precise and concise.

But very different views of the nature of intelligent reasoning are also possible. One distinctly different view is embedded in the part of AI that is influenced by the psychological tradition. This tradition, rooted in the work of D. O. Hebb, J. Bruner, G. Miller, and A. Newell and H. Simon, broke through the stimulus-response view demanded by behaviorism and suggested instead that human problem-solving behavior could usefully be viewed in terms of goals, plans, and other complex mental structures. Modern manifestations include work on SOAR as a general mechanism for producing intelligent reasoning and knowledge-based systems as a means of capturing human expert reasoning.

Comparing these two traditions reveals significant differences and illustrates the consequences of adopting one or the other view of intelligent reasoning. In the logicist tradition intelligent reasoning is taken to be a form of calculation, typically, deduction in first-order logic, while the tradition based in

psychology takes as the defining characteristic of intelligent reasoning that it is a particular variety of human behavior. In the logicist view, the object of interest is, thus, a construct definable in formal terms through mathematics, while for those influenced by the psychological tradition, it is an empirical phenomenon from the natural world. Thus, there are two very different assumptions here about the essential nature of the fundamental phenomenon to be captured.

A second contrast arises in considering the character of the answers each seeks. The logicist view has traditionally sought compact and precise characterizations of intelligence, looking for the kind of characterizations encountered in mathematics (and at times in physics). By contrast, the psychological tradition suggests that intelligence is not only a natural phenomenon, it is also an inherently complex natural phenomenon: As human anatomy and physiology are inherently complex systems resulting from a long process of evolution, so perhaps is intelligence. As such, intelligence may be a large and fundamentally ad hoc collection of mechanisms and phenomena, one that complete and concise descriptions might not be possible for.

Several useful consequences result from understanding the different positions on this fundamental question that are taken by each tradition. First, it demonstrates that selecting any of the modern offspring of these traditions—that is, any of the representation technologies shown at the bottom of the table—means choosing more than a representation. In the same act, we are also selecting a conception of the fundamental nature of intelligent reasoning.

Second, these conceptions differ in important ways: There are fundamental differences in the conception of the phenomenon we are trying to capture. The different conceptions in turn mean there are deep-seated differences in the character and the goals of the various research efforts that are trying to create intelligent programs. Simply put, different conceptions of the nature of intelligent reasoning lead to different goals, definitions of success, and different artifacts being created.

Finally, these differences are rarely articulated. In turn, this lack of articulation leads to arguments that may be phrased in terms of issues such as representation choice (for

The choice of appropriate vocabulary and the degree of formality depends, in turn, on the basic conception of intelligent behavior

example, the virtues of sound reasoning in first-order predicate calculus versus the difficult-to-characterize inferences produced by frame-based systems) when the real issues are, we believe, the different conceptions of the fundamental nature of intelligence. Understanding the different positions assists in analyzing and sorting out the issues appropriately.

Which Inferences Are Sanctioned? The second component of a representation's theory of intelligent reasoning is its set of sanctioned inferences, that is, a selected set of inferences that are deemed appropriate conclusions to draw from the information available. The classic definition is supplied by traditional formal logic, where the only sanctioned inferences are sound inferences (those encompassed by logical entailment, in which every model for the axiom set is also a model for the conclusion). This answer has a number of important benefits, including being intuitively satisfying (a sound argument never introduces error), explicit (so we know precisely what we're talking about), precise enough that it can be the subject of formal proofs, and old enough that we have accumulated a significant body of experience with it.

Logic has also explored several varieties of unsound inference, including circumscription and abduction. This exploration has typically been guided by the requirement that there be "a well motivated model-theoretic justification" (Nilsson 1991, pp. 42–43), such as the minimal model criterion of circumscription. This requirement maintains a fundamental component of the logicist approach: Although it is willing to arrive at conclusions that are true in some subset of the models (rather than true in every model), the set of sanctioned inferences is still conceived of in model-theoretic terms and is specified precisely in these terms.

Other representations have explored other definitions: probabilistic reasoning systems (for example, Pearl [1988]) sanction the inferences specified by probability theory, while work on rational agents (for example, Doyle [1992]) relies on concepts from the theory of economic rationality.

Among the common knowledge representation technologies, rule-based systems capture guesses of the sort that a human expert makes, guesses that are not necessarily either sound or true in any model. A frame-based representation encourages jumping to possibly incorrect conclusions based on good matches, expectations, or defaults. Both of

these representations share the psychological tradition of defining the set of sanctioned inferences with reference to the behavior of the human expert rather than reference to an abstract formal model.

As these examples show, different approaches to representation specify sanctioned inferences in ways that differ in both content and form. Where the specification for logic, for example, is expressed in terms of model theory and is mathematically precise, other representations provide answers phrased in other terms, often with considerably less precision. Frames theory, for example, offers a definition phrased in terms of human behavior and is specified only approximately.

The differences in both content and style in turn have their origin in the different conceptions of intelligent reasoning that were explored previously. Phrasing the definition in terms of human behavior is appropriate for frames because the theory conceives of intelligent reasoning as a characteristic form of human behavior. In attempting to describe this behavior, the theory is faced with the task of characterizing a complex empirical phenomenon that can be captured only roughly at the moment and that might never be specifiable with mathematical precision, hence the appropriateness of an approximate answer.

For frames theory then, the specification of sanctioned inferences is both informal and empirical, as an unavoidable consequence of its conception of intelligence. The work (and other work like it) is neither sloppy nor causally lacking in precision; the underlying conception of intelligent reasoning dictates a different approach to the task, a different set of terms in which to express the answer, and a different focus for the answer.

The broader point here is to acknowledge the legitimacy of a variety of approaches to specifying sanctioned inferences: Model theory might be familiar and powerful, but even for formal systems, it is not the only possible language. More broadly still, formal definitions are not the only terms in which the answer can be specified. The choice of appropriate vocabulary and the degree of formality depends, in turn, on the basic conception of intelligent behavior.

Which Inferences Are Recommended? While sanctioned inferences tell us what conclusions we are permitted to make, this set is invariably very large and, hence, provides insufficient constraint. Any automated system attempting to reason, guided only by

knowing what inferences are sanctioned, soon finds itself overwhelmed by choices. Hence, we need more than an indication of which inferences we can legally make; we also need some indication of which inferences are appropriate to make, that is, intelligent. This indication is supplied by the set of recommended inferences.

Note that the need for a specification of recommended inferences means that in specifying a representation, we also need to say something about how to reason intelligently. Representation and reasoning are inextricably and usefully intertwined: A knowledge representation is a theory of intelligent reasoning.

This theory often results from observation of human behavior. Minsky's original exposition of frame theory, for example, offers a clear example of a set of recommended inferences inspired by observing human behavior. Consider the following statement from Minsky's abstract (1974, 1975) to his original frames paper:

This is a partial theory of thinking.... Whenever one encounters a new situation (or makes a substantial change in one's viewpoint), he selects from memory a structure called a *frame*; a remembered framework to be adapted to fit reality by changing details as necessary.

A frame ... [represents] a stereotyped situation, like being in a certain kind of living room, or going to a child's birthday party.

The first sentence illustrates the intertwining of reasoning and representation: This paper is about knowledge representation, but it announces at the outset that it is also a theory of thinking. In turn, this theory arose from an insight about human intelligent reasoning, namely, how people might manage to make the sort of simple commonsense inferences that appear difficult to capture in programs. The theory singles out a particular set of inferences to recommend, namely, reasoning in the style of anticipatory matching.

Similar characterizations of recommended inferences can be given for most other representation technologies. Semantic nets in their original form, for example, recommend bidirectional propagation through the net, inspired by the interconnected character of word definitions and the part of human intelligence manifested in the ability of people to find connections between apparently disparate concepts. The rules in knowledge-based systems recommend plausible inferences, inspired by the observation of

human expert reasoning.

By contrast, logic has traditionally taken a minimalist stance on this issue. The representation itself offers only a theory of sanctioned inferences, seeking to remain silent on the question of which inferences to recommend.

The silence on this issue is motivated by a desire for generality in the inference machinery and a declarative (that is, use-dependent) form for the language, both fundamental goals of the logicist approach: "... logicists strive to make the inference process as uniform and domain independent as possible and to represent all knowledge (even the knowledge about how to use knowledge) declaratively" (Nilsson 1991, p. 46).

But a representation with these goals cannot single out any particular set of inferences to recommend for two reasons. First, if the inference process is to be general and uniform (that is, work on all problems and work in the same way), it must be neutral about which inferences to recommend; any particular subset of inferences it attempted to single out might be appropriate in one situation but fatally bad in another because no inference strategy (unit preference, set of support, and so on) is universally appropriate. Second, if statements in the language are to be declarative, they must express a fact without any indication of how to reason with it (use-free expression is a defining characteristic of a declarative representation). Hence, the inference engine can't recommend any inferences (or it loses its generality and uniformity), and the statements of fact in the language cannot recommend any inferences (because by embedding such information, they lose their declarative character).⁴

Thus, the desire for generality and use-free expression prevents the representation itself from selecting inferences to recommend. But if the representation itself cannot make the recommendation, the user must because the alternative—unguided search—is untenable.

Requiring the user to select inferences is, in part, a deliberate virtue of the logicist approach: Preventing the representation from selecting inferences and, hence, requiring the user to do so offers the opportunity for this information to be represented explicitly rather than embedded implicitly in the machinery of the representation (as, for example, in rule-based systems or PROLOG).

One difficulty with this admirable goal arises in trying to provide the user with the tools to express the strategies and guide the system. Three approaches are commonly used: (1) have the user tell the system what to

the desire for generality and use-free expression prevents the representation itself from selecting inferences to recommend

do, (2) have the user lead it into doing the right thing, and (3) build in special-purpose inference strategies. By *telling the system what to do*, we mean that the user must recommend a set of inferences by writing statements in the same (declarative) language used to express facts about the world (for example, MRS [Russell 1985]). By *leading the system into doing the right thing*, we mean that the user must carefully select the axioms, theorems, and lemmas supplied to the system. The presence of a lemma, for example, is not simply a fact the system should know; it also provides a way of abbreviating a long chain of deductions into a single step, in effect allowing the system to take a large step in a certain direction (namely, the direction in which the lemma takes us). By carefully selecting facts and lemmas, the user can indirectly recommend a particular set of inferences. By *special-purpose inference strategies*, we mean building specific control strategies directly into the theorem prover. This

purposely silent on the issue of recommended inferences, logic offers both a degree of generality and the possibility of making information about recommended inferences explicit and available to be reasoned about in turn. On the negative side, the task of guiding the system is left to the user, with no conceptual assistance offered, and the practices that result at times defeat some of the key goals that motivated the approach at the outset.

Role 4: A Knowledge Representation Is a Medium for Efficient Computation

From a purely mechanistic view, reasoning in machines (and, perhaps, in people) is a computational process. Simply put, to use a representation, we must compute with it. As a result, questions about computational efficiency are inevitably central to the notion of representation.

This fact has long been recognized, at least implicitly, by representation designers: Along



The good news here is that by remaining purposely silent on the issue of recommended inferences, logic offers both a degree of generality and the possibility of making information about recommended inferences explicit and available to be reasoned about in turn

approach can offer significant speedup and a pragmatically useful level of computational efficiency.

Each of these approaches has both benefits and drawbacks. Expressing reasoning strategies in first-order logic is in keeping with the spirit of the logicist approach, namely, explicit representation of knowledge in a uniform, declarative representation. But this approach is often problematic in practice: a language designed to express facts declaratively is not necessarily good for expressing the imperative information characteristic of a reasoning strategy.

Careful selection of lemmas is, at best, an indirect encoding of the guidance information to be supplied. Finally, special-purpose deduction mechanisms are powerful but embed the reasoning strategy both invisibly and procedurally, defeating the original goals of domain-independent inference and explicit, declarative representation.

The good news here is that by remaining

with their specification of a set of recommended inferences, representations typically offer a set of ideas about how to organize information in ways that facilitate making these inferences. A substantial part of the original frames notion, for example, is concerned with just this sort of advice, as more of the frames paper illustrates (Minsky 1974, 1975):

A frame ... [represents] a stereotyped situation, like being in a certain kind of living room, or going to a child's birthday party.

Attached to each frame are several kinds of information. Some of this information is about how to use the frame. Some is about what one can expect to happen next. Some is about what to do if these expectations are not confirmed.

The notion of triggers and procedural attachment in frames is not so much a statement about what procedures to write (the

theory is rather vague here) as it is a description of a useful way to organize information, for example (paraphrasing the previous quotation), attach to each frame information about how to use the frame and what to do if expectations are not confirmed. Similarly, organizing frames into taxonomic hierarchies both suggests taxonomic reasoning and facilitates its execution (as in structured inheritance networks).

Other representations provide similar guidance. Traditional semantic nets facilitate bi-directional propagation by the simple expedient of providing an appropriate set of links, while rule-based systems facilitate plausible inferences by supplying indexes from goals to rules whose conclusion matches (backward chaining) and from facts to rules whose premise matches (forward chaining).

While the issue of efficient use of representations has been addressed by representation designers, in the larger sense, the field appears to have been historically ambivalent in its reaction. Early recognition of the notion of heuristic adequacy (McCarthy and Hayes 1969) demonstrates that early on, researchers appreciated the significance of the computational properties of a representation, but the tone of much subsequent work in logic (for example, Hayes [1979]) suggested that *epistemology* (knowledge content) alone mattered and defined computational efficiency out of the agenda. Of course, epistemology does matter, and it can be useful to study it without the potentially distracting concerns about speed. But eventually, we must compute with our representations; hence efficiency must be part of the agenda.

The pendulum later swung sharply over to what we might call the computational imperative view. Some work in this vein (for example, Levesque and Brachman [1985]) offered representation languages whose design was strongly driven by the desire to provide not only efficiency but also guaranteed efficiency. The result appears to be a language of significant speed but restricted power (Doyle 1991, 1989).

Either end of this spectrum seems problematic: We ignore computational considerations at our peril, but we can also be overly concerned with them, producing representations that are fast but inadequate for real use.

Role 5: A Knowledge Representation Is a Medium of Human Expression

Finally, knowledge representations are also the means by which we express things about the world, the medium of expression and

communication in which we tell the machine (and perhaps one another) about the world. This role for representations is inevitable as long as we need to tell the machine (or other people) about the world and as long as we do so by creating and communicating representations.⁵ Thus, the fifth role for knowledge representations is as a medium of expression and communication for our use.

In turn, this role presents two important sets of questions. One set is familiar: How well does the representation function as a medium of expression? How general is it? How precise? Does it provide expressive adequacy? and so on.

An important question that is discussed less often is, How well does it function as a medium of communication? That is, how easy is it for us to talk or think in this language? What kinds of things are easily said in the language, and what kinds of things are so difficult that they are pragmatically impossible?

Note that the questions here are of the form, How easy is it? rather than, Can we? This language is one that we must use; so, things that are possible in principle are useful but insufficient; the real question is one of pragmatic utility. If the representation makes things possible but not easy, then as real users we might never know whether we misunderstood the representation and just do not know how to use it or whether it truly cannot express some things that we would like to say. A representation is the language in which we communicate; hence, we must be able to speak it without heroic effort.

Consequences for Research and Practice

We believe that this view of knowledge representation can usefully influence practice and can help inform the debate surrounding several issues in representation research. For practice, it offers a framework that aids in making explicit the important insights and spirit of a representation and illustrates the difference in design that results from indulging, rather than violating, this spirit.

The consequences of the view for research include (1) a broader conception of representation, urging that all the roles should be kept in mind when creating representation languages, (2) the recognition that a representation embeds a theory of intelligent reasoning, (3) the ability to use the broader view of representation to guide the combination of representations, (4) the ability to use the broader

view to dissect some of the arguments about formal equivalence of representations, and (5) the belief that the central task of knowledge representation is capturing the complexity of the real world.

Space limitations require that we only briefly sketch out these consequences here. A complete discussion is found in Davis, Shrobe, and Szolovits (1993).

Consequence for Practice: Characterizing the Spirit of a Representation

The roles enumerated previously help to characterize and make explicit the spirit of a representation, that is, the important set of ideas and inspirations that lie behind (and, significantly, are often less obvious than) the concrete machinery used to implement the representation. The spirit is often difficult to describe with precision, but we believe it is well characterized by the last four of the roles we just enumerated (all representations are surrogates; hence, there is little difference among them on the first role). The stance that a representation takes on each of these issues, along with its rationale for this stance, indicates what the representation is trying to say about how to view and reason about the world.

In its original incarnation (Minsky 1974, 1975), the frames idea, for example, is primarily an ontological commitment and a theory of intelligent reasoning based on insights about human cognition and the organization of knowledge in memory. The major ontological commitment is to view the world in terms of *stereotypical descriptions*, that is, concepts described in terms of what is typically true about them. This approach is particularly well suited to concepts in the natural world, where categories rarely have precise specifications in terms of necessary and sufficient conditions, and exceptions abound. There is additional ontological commitment in linking frames into systems to capture perspective shifts: We are encouraged to look for such shifts when viewing the world.

The theory of intelligent reasoning embedded in frames claims that much reasoning is based on recognition, particularly matching stereotypes against individual instances. The suggestions concerning the organization of knowledge are based on the belief that information in human memory is richly and explicitly interconnected rather than structured as a set of independent or only implicitly connected facts. Thus, the frames theory

recommends inferences produced by stereotype matching and instantiation and facilitates these inferences through the frame structure itself as well as the further organization of frames into frame systems.

The theory sanctions inferences that are unsound, as in the analogical and default reasoning done when matching frames. It also sanctions inferences that involve relatively large mismatches in order to model understanding that is tenacious even in the face of inconsistencies.

The theory provides a medium for potentially efficient computation by casting understanding as matching rather than deduction. Finally, it offers a medium of expression that is particularly useful for describing concepts in the natural world, where we often need some way of indicating what properties an object typically has, without committing to statements about what is always true.

Two useful consequences result from characterizing a representation in these terms, making its position on the roles both explicit and understandable: first, it enables a kind of explicitly invoked Whorfian theory of representation use. Although the representation we select will have inevitable consequences for how we see and reason about the world, we can at least select it consciously and carefully, trying to find a pair of glasses appropriate for the task at hand. Steps in this direction include having representation designers carefully characterize the nature of the glasses they are supplying (for example, making explicit the ontological commitments, recommended inferences) and having the field develop principles for matching representations to tasks.

Second, such characterizations would facilitate the appropriate use of a representation. By *appropriate*, we mean using it in its intended spirit, that is, using it for what it was intended to do, not for what it can be made to do. Yet with striking regularity, the original spirit of a representation is seen as an opponent to be overcome. With striking regularity, the spirit is forgotten, replaced by a far more mechanistic view that sees a data structure rather than a representation, computation rather than inference. Papers written in this mind set typically contain claims of how the author was able, through a creative, heroic, and often obscure act, to get a representation to do something we wouldn't ordinarily have thought it capable of doing.

However, if such obscure acts are what

using a representation is all about, it becomes an odd and terribly awkward form of programming: The task becomes one of doing what we already know we want to do but being forced to do it using just the constructs available in the representation, no matter how good or bad the fit.

In this case, the creative work is often in overcoming the representation, seeing how we can get it to behave like something else.⁶ The result is knowledge representations applied in ways that are uninformed by the inspirations and insights that led to their invention and that are the source of their power. Systems built in this spirit often work despite their representations, not because of them; they work because the authors, through great effort, managed to overcome the representation.

Consequence for Research: Representation and Reasoning Are Intertwined

At various times in the development of the field, the suggestion has been made that we ought to view knowledge representation in purely epistemological terms; that is, take the singular role of representation to be conveying knowledge content (for example, Hayes [1979]). As we noted earlier, epistemology matters, but it is not the whole of the matter. Representation and reasoning are inextricably intertwined: We cannot talk about one without also unavoidably discussing the other. We argue as well that the attempt to deal with representation as knowledge content alone leads to an incomplete conception of the task of building an intelligent reasoner.

Each of these claims is grounded in an observation made earlier. We observed first that every representation embeds at its core a conception of what constitutes intelligent reasoning (table 1). Hence, any discussion of representation unavoidably carries along with it a (perhaps implicit) view of intelligent reasoning.

We also observed that in building an intelligent reasoner, it is not enough to indicate what inferences are legal; we also need to know which are appropriate (that is, recommended). A familiar example from logic makes the point nicely: From A , we can infer $A \wedge A$, $A \wedge A \wedge A$, and so on. All these inferences are legal, but they are hardly intelligent.

Hence, we arrive at our claim that a theory of legal (sanctioned) inference is insufficient; to build an intelligent reasoner, we also need a theory of intelligent inference. In fact, there

might be multiple theories of intelligent inference, each specific to a particular task domain.

Consequence for Research: Combining Representations

Recall that a representation is, among other things, a theory of intelligent reasoning and a collection of mechanisms for implementing this theory. We believe that appropriate attention to both of these aspects, in the appropriate order, makes a significant difference in the outcome of any effort at representation combination.

Too often, efforts at combination appear to be conceived of in terms of finding ways for the two mechanisms to work together, with insufficient (and sometimes no) attention to what we consider to be the much more important task: determining how the two theories of intelligent reasoning might work together. Focusing on mechanisms means determining such things as how, say, rules, procedures, and objects might invoke one another interchangeably. Focusing on the theories of intelligent reasoning means attending to what kinds of reasoning are within the spirit of each representation and how these varieties of reasoning might sensibly be combined.

Two efforts at combining representations illustrate the different conceptions of the task that arise from focusing on mechanism and focusing on reasoning. As the first example, consider this description of the *LOOPS* system (Stefik et al. 1983) and its efforts at integrating several paradigms:

Some examples illustrate the integration of paradigms in *LOOPS*: the “workspace” of a ruleset is an object, rules are objects, and so are rulesets. Methods in classes can be either Lisp functions or rulesets. The procedures in active values can be Lisp functions, rulesets, or calls on methods. The ring in the *LOOPS* logo reflects the fact that *LOOPS* not only contains the different paradigms, but integrates them. The paradigms are designed not only to compliment each other, but also to be used together in combination. (p. 5)

Contrast the mind set and the previous approach with this view of a similar undertaking, also aimed at combining rules and frames (Yen, Neches, and MacGregor 1989).

Rules and frames are two contrasting schemes for representing different kinds of knowledge. Rules are appropriate for representing logical implications, or for associating actions with conditions under which the actions should be

There is a familiar pattern in knowledge representation research in which the description of a new knowledge representation technology is followed by claims that the new ideas are, in fact, formally equivalent to an existing technology.

taken.... Frames (or semantic nets) are appropriate for defining terms and for describing objects and the taxonomic class/membership relationships among them. An important reasoning capability of frame systems with well-defined semantics is that they can infer the class/membership relations between frames based on their definitions.

Since the strengths and weaknesses of rules and frames are complementary to each other, a system that integrates the two will benefit from the advantages of both techniques. This paper describes a hybrid architecture called *classification based programming* which extends the production system architecture using automatic classification capabilities within frame representations. In doing so, the system enhances the power of a pattern matcher in a production system from *symbolic matching* to *semantic matching*, organizes rules into rule classes based on their functionalities, and infers the various relationships among rules that facilitate explicit representation of control knowledge. (p. 2)

Note, in particular, how the first of these efforts focuses on computational mechanisms, while the second is concerned with representation and reasoning. The first seeks to allow several programming constructs—among them, rules and structured objects—to work together, while the second attempts to allow two representations—rules and structured objects—to work together. The first seeks to permit mechanisms to invoke one another, while the second considers the different varieties of inference natural to two representations and suggests how these two kinds of reasoning could work synergistically: Rules are to be used for the kind of reasoning they capture best—unrestricted logical implications—and frames are to be used for their strength, namely, taxonomic reasoning. Thus, the first paper (Stefik et al. 1983) proposes a computational architecture, while the second (Yen, Neches, and MacGregor 1989) offers a representation and reasoning architecture.

Both of these undertakings are important, but we believe that where the goal is combining representations, the task should be conceived of in terms central to a representation: its theory of intelligent reasoning. To do this, we should consider what kind of reasoning we expect from each representation, we should propose a design for how these reasoning schemes will work together, and we

should have a rationale for believing that this combination will be effective. The first and most important thing we require in this task is a representation architecture; from that the appropriate computational architecture follows. We believe that efforts that attend only to achieving the appropriate computational architecture are beginning the effort in the wrong place and will fall far short of a crucial part of the goal.

Consequence for Research: Arguments about Formal Equivalence

There is a familiar pattern in knowledge representation research in which the description of a new knowledge representation technology is followed by claims that the new ideas are, in fact, formally equivalent to an existing technology. Historically, the claim has often been phrased in terms of equivalence to logic. Semantic nets, for example, have been described in these terms (Hayes 1977; Nash-Webber and Reiter 1977), while the development of frames led to a sequence of such claims, including the suggestion that “most of ‘frames’ is just a new syntax for parts of first-order logic” (Hayes 1979, p. 293).

That frames might also be an alternative syntax seems clear; that they are merely or just an alternative syntax seems considerably less clear. That frames are not entirely distinct from logic seems clear; that all of the idea can be seen as logic seems considerably less clear.

We believe that claims such as these are substantive only in the context of a narrowly construed notion of what a knowledge representation is. Hayes (1979) is explicit about part of his position on a representation: “One can characterise a representational language as one which has (or can be given) a semantic theory” (p. 288). He is also explicit about the tight lines drawn around the argument: “Although frames are sometimes understood at the metaphysical level and sometimes at the computational level, I will discuss them as a representational proposal” (p. 288), that is, as a language with a semantic theory and nothing more. Both metaphysics and computation have been defined as out of the agenda. Hayes also says, “None of this discussion [about frames as a computational device for organizing memory access and inference] engages representational issues” (p. 288). Here, it becomes evident that the claim is less about frames and more a definition of what will be taken as representational issues.

Note that specifically excluded from the discussion are the ontological commitment

of the representation, namely, “what entities shall be assumed to exist in the world” (Hayes 1979, p. 288), and the computational properties that the representation provides. Despite claims to the contrary, we argue that the ontology of frames (and other representations) and computational questions not only engage representational issues, they are representational issues. These and other properties are crucial to knowledge representation both in principle and in any real use.

Consequences for Research: All Five Roles Matter

Although representations are often designed with considerable attention to one or another of the issues listed in our five roles, we believe that all the roles are of substantial significance and that ignoring any one of them can lead to important inadequacies. While designers rarely overlook representation’s role as a surrogate and as a definition of sanctioned inferences, insufficient guidance on the other roles is not uncommon, and the consequences matter.

As we argued earlier, for example, ontological commitment matters: The guidance it provides in making sense of the profusion of detail in the world is among the most important things a representation can supply. Yet some representations, in a quest for generality, offer little support on this dimension.

A similar argument applies to the theory of intelligent reasoning: A representation can guide and facilitate reasoning if it has at its heart a theory of what reasoning to do. Insufficient guidance here leaves us susceptible to the traditional difficulties of unguided choice.

Pragmatically efficient computation matters because most of the use of a representation is (by definition) in the average case. Interest in producing weaker representations to guarantee improved worst-case performance may be misguided, demanding far more than is necessary and paying a heavy price for it.⁷

The use of a representation as a medium of expression and communication matters because we must be able to speak the language to use it. If we can’t determine how to say what we’re thinking, we can’t use the representation to communicate with the reasoning system.

Attempting to design representations to accommodate all five roles is, of course, challenging, but we believe the alternative is the creation of tools with significant deficiencies.

The Goal of Knowledge Representation Research

We believe that the driving preoccupation of the field of knowledge representation should be understanding and describing the richness of the world. Yet in practice, research that describes itself as core knowledge representation work has concentrated nearly all its efforts in a much narrower channel, much of it centered around taxonomic and default reasoning (for example, Brachman and Schmolze [1985]; Levesque and Brachman [1985]; Fahlman, Touretsky, and van Roggen [1981]).

We believe that it is not an accident that a useful insight about finding a good set of temporal abstractions came from close examination of a realistic task set in a real-world domain (Hamscher 1991). It underscores our conviction (shared by others; see Lenat [1990]) that attempting to describe the richness of the natural world is the appropriate forcing function for knowledge representation work.

Our point here concerns both labeling and methodology: (1) work such as Hamscher (1991) and Lenat (1990) should be recognized by the knowledge representation community as of central relevance to knowledge representation research, not categorized as diagnosis or qualitative physics and seen as unrelated, and (2) insights of the sort obtained in Hamscher (1991) and Lenat (1990) come from studying the world, not from studying languages. We argue that those who choose to identify themselves as knowledge representation researchers should be developing theory and technology that facilitate projects such as these, and conversely, those who are building projects such as these are engaged in a centrally important variety of knowledge representation research.

While tools and techniques are important, the field is and ought to be much richer than that, primarily because the world is much richer than that. We believe that understanding and describing this richness should be the central preoccupation of the field.

Summary

We argued that a knowledge representation plays five distinct roles, each important to the nature of representation and its basic tasks. These roles create multiple, sometimes competing demands, requiring selective and intelligent trade-offs among the desired characteristics. These five roles also aid in clearly characterizing the spirit of the repre-

*the
fundamental
task of
representation
is describing
the natural
world ...*

sentations and the representation technologies that have been developed.

This view has consequences for both research and practice in the field. On the research front, it argues for a conception of representation that is broader than the one often used, urging that all five aspects are essential representation issues. It argues that the ontological commitment that a representation supplies is one of its most significant contributions; hence, the commitment should be both substantial and carefully chosen. It also suggests that the fundamental task of representation is describing the natural world and claims that the field would advance furthest by taking this view as its central preoccupation.

For the practice of knowledge representation work, the view suggests that combining representations is a task that should be driven by insights about how to combine their theories of intelligent reasoning, not their implementation mechanisms. The view also urges the understanding of and indulgence of the fundamental spirit of representations. We suggest that representation technologies should not be considered as opponents to be overcome, forced to behave in a particular way, but instead, they should be understood on their own terms and used in ways that rely on the insights that were their original inspiration and source of power.

Acknowledgments

This article describes research done at the Artificial Intelligence Laboratory and the Laboratory for Computer Science at the Massachusetts Institute of Technology (MIT). Support for this work was received from the Defense Advanced Research Projects Agency under Office of Naval Research contract N00014-85-K-0124; the National Library of Medicine under grant R01 LM 04493; the National Heart, Lung, and Blood Institute under grant R01 HL 33041; Digital Equipment Corporation; and the General Dynamics Corp.

Earlier versions of this material formed the substance of talks delivered at the Stanford University Knowledge Systems Laboratory and an invited talk given at the Ninth National Conference on AI; many useful comments from the audience were received on both occasions. Jon Doyle and Ramesh Patil offered a number of insightful suggestions on drafts of this article; we are also grateful to Rich Fikes, Pat Hayes, Ron Brach-

man, and Hector Levesque for extended discussions about some of this material.

Notes

1. Conversely, action can substitute for reasoning. This dualism offers one way of understanding the relation between traditional symbolic representations and the *situated-action approach*, which argues that action can be linked directly to perception, without the need for intermediating symbolic representations.

2. The phrase ontological commitment is perhaps not precisely correct for what we have in mind here, but it is the closest available approximation. While ontology is, strictly speaking, concerned with what exists in the world, we phrased this section carefully in terms of how to view the world, purposely sidestepping many standard thorny philosophical issues surrounding claims of what exists. A second way around the issue is to note that the world we are interested in capturing is the world inside the mind of some intelligent human observer (for example, a physician, an engineer); in which case, it can plausibly be argued that in this world, rules, prototypes, and so on, do exist.

3. Note that even at the outset, there is a hint that the desired form of reasoning might be describable in a set of formal rules.

4. The consequences of this approach are evident even in the use of disjunctive normal form as a canonical representation: Although

$$X1 \wedge X2 \wedge X3 \rightarrow X4$$

is semantically equivalent to

$$\neg X1 \vee \neg X2 \vee X4 \vee \neg X3,$$

some potentially useful information is lost in the transformation. The first form might suggest that $X1$, $X2$, and $X3$ have something in common, namely, that they should be thought of as the preconditions needed to establish $X4$. This hint might be useful in deciding how to reason in the problem, but if so, it is lost in the transformation to disjunctive normal form. By contrast, consider languages such as MICROPLANNER and PROLOG, which make explicit use of the form of the inference rule to help guide the deduction process.

5. It will presumably continue to be useful even if machines invent their own knowledge representations based on independent experience of the world. If their representations become incomprehensible to us, the machines will be unable to either tell us what they know or explain their conclusions.

6. Of course, there is utility in establishing the equivalence of two representations by showing how one can be made to behave like another. But this exercise needs to be done only once, and it is done for its own sake rather than because it is good practice in system construction.

7. As we argued elsewhere (Davis 1991), a *computational cliff* (that is, unacceptable worst-case behavior) is a problem only if every inference, once set in motion, cannot possibly be interrupted. The simple expedient of resource-limited computation prevents any inference from permanently trapping the pro-

gram in a task requiring unreasonable amounts of time.

References

- Brachman, R., and Levesque, H. eds. 1985. *Readings in Knowledge Representation*. San Mateo, Calif.: Morgan Kaufmann.
- Brachman, R., and Schmolze, J. 1985. An Overview of the KL-ONE Knowledge Representation System. *Cognitive Science* 9(2): 171–216.
- Davis, R. 1991. A Tale of Two Knowledge Servers. *AI Magazine* 12(3): 118–120.
- Davis, R., and Shrobe, H. 1983. Representing Structure and Behavior of Digital Hardware. *IEEE Computer* (Special Issue on Knowledge Representation) 16(10): 75–82.
- Davis, R.; Shrobe, H.; and Szolovits, P. 1993. What Is a Knowledge Representation? Memo, AI Laboratory, Massachusetts Institute of Technology.
- Doyle, J. 1992. Rationality and Its Roles in Reasoning. *Computational Intelligence* 8:376–409.
- Doyle, J., and Patil, R. 1991. Two Theses of Knowledge Representation: Language Restrictions, Taxonomic Classification, and the Utility of Representation Services. *Artificial Intelligence* 48(3): 261–297.
- Doyle, J., and Patil, R. 1989. Two Dogmas of Knowledge Representation, Technical Memo, 387B, Laboratory for Computer Science, Massachusetts Institute of Technology.
- Fahlman, S.; Touretsky, D.; and van Roggen, W. 1981. Cancellation in a Parallel Semantic Network. In *Proceedings of the Seventh International Joint Conference on Artificial Intelligence*, 257–263. Menlo Park, Calif.: International Joint Conferences on Artificial Intelligence.
- Hamscher, W. 1991. Modeling Digital Circuits for Troubleshooting. *Artificial Intelligence* 51:223–272.
- Hayes, P. 1979. The Logic of Frames. In *Readings in Knowledge Representation*, eds. R. Brachman and H. Levesque, 288–295. San Mateo, Calif.: Morgan Kaufmann.
- Hayes, P. 1978. Naive Physics I: Ontology for Liquids. In *Formal Theories of the Commonsense World*, eds. J. R. Hobbs and R. C. Moore. Norwood, N.J.: Ablex.
- Hayes, P. 1977. In Defense of Logic. In *Proceedings of the fifth International Joint Conference on Artificial Intelligence*, 559–565. Menlo Park, Calif.: International Joint Conferences on Artificial Intelligence.
- Lenat, D.; Guha, R.; Pittman, K.; Pratt, D.; and Shepherd, M. 1990. cyc: Toward Programs with Common Sense. *Communications of the ACM* 33(8): 30–49.
- Levesque, H., and Brachman, R. 1985. A Fundamental Trade-Off in Knowledge Representation and Reasoning. In *Readings in Knowledge Representation*, eds. R. Brachman and H. Levesque, 42–70. San Mateo, Calif.: Morgan Kaufmann.
- McCarthy, J., and Hayes, P. 1969. Some Philosophical Problems from the Standpoint of AI. In *Machine Intelligence 4*, eds. B. Meltzer and D. Michie, 463–504. Edinburgh: Edinburgh University Press.
- Minsky, M. 1975. A Framework for Representing Knowledge. In *The Psychology of Computer Vision*, ed. P. Winston, 211–277. New York: McGraw-Hill.
- Minsky, M. 1974. A Framework for Representing Knowledge, Memorandum, 306, AI Laboratory, Massachusetts Institute of Technology.
- Nash-Webber, B., and Reiter, R. 1977. Anaphora and Logical Form. In *Proceedings of the fifth International Joint Conference on Artificial Intelligence*, 121–131. Menlo Park, Calif.: International Joint Conferences on Artificial Intelligence.
- Nilsson, N. 1991. Logic and Artificial Intelligence. *Artificial Intelligence* 47(1): 31–56.
- Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Mateo, Calif.: Morgan Kaufmann.
- Pople, H. 1982. Heuristic Methods for Imposing Structure on Ill-Structured Problems. In *AI in Medicine*, ed. P. Szolovits, 119–190. American Association for the Advancement of Science Symposium 51. Boulder, Colo.: Westview.
- Russell, S. 1985. The Compleat Guide to MRS, Technical Report, STAN-CS-85-1080, Dept. of Computer Science, Stanford Univ.
- Stefik, M.; Bobrow, D.; Mittal, S.; and Conway, L. 1983. Knowledge Programming in LOOPS: Report on an Experimental Course. *AI Magazine* 4(3): 3–13.
- Yen, J.; Neches, R.; and MacGregor, R. 1989. Classification-Based Programming: A Deep Integration of Frames and Rules, Report ISI/RR-88-213, USC/Information Sciences Institute, Marina del Rey, California.

Randall Davis is a professor in the Electrical Engineering and Computer Science Department at the Massachusetts Institute of Technology (MIT) and associate director of the MIT AI Laboratory. He and his group at MIT have produced a variety of model-based reasoning systems for electric and (more recently) mechanical systems.

Howard Shrobe is a principle research scientist at the Massachusetts Institute of Technology (MIT) AI Laboratory and a technical director at Symbolics Inc. He has conducted research at MIT on the use of knowledge-based systems in engineering and design. Through Symbolics, he has participated in the implementation and deployment of several large-scale expert systems. He is a fellow of the American Association for Artificial Intelligence.

Peter Szolovits is associate professor of computer science and engineering at the Massachusetts Institute of Technology (MIT) and head of the Clinical Decision-Making Group within the MIT Laboratory for Computer Science. Szolovits's research centers on the application of AI methods to problems of medical decision making. He has worked on problems of diagnosis of kidney diseases, therapy planning, execution and monitoring for various medical conditions, and computational aspects of genetic counseling. His interests in AI include knowledge representation, qualitative reasoning, and probabilistic inference.