

**NOTICE WARNING CONCERNING COPYRIGHT RESTRICTIONS:**  
The copyright law of the United States (title 17, U.S. Code) governs the making of photocopies or other reproductions of copyrighted material. Any copying of this document without permission of its author may be prohibited by law.

University Libraries  
Carnegie Mellon University  
Pittsburgh PA 15213-3890

This research was supported by the Defense Advanced Research Projects Agency (DOD), ARPA Order No. 3597, monitored by the Air Force Avionics Laboratory under Contract F33615-81-K-1539. The views and conclusions contained in this document are those of the author, and should not be interpreted as representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the Government of the United States of America.

# The Foundations of Psychology

A logico-computational inquiry into the concept of mind

by Dr. Jon Doyle  
Research Associate at Carnegie-Mellon University

October 1981-February 1982  
Pittsburgh, Pennsylvania  
© Copyright 1981, 1982 by Jon Doyle

February 18, 1982

## Abstract

We compare certain trains of thought in philosophy of mind and artificial intelligence which lead to a remarkable convergence of ideas. Demands from philosophy that psychological theories have predictive power join with demands from artificial intelligence that machines adaptively maintain their own mental state to suggest a conception of minds as narrowly realized psychological theories. We use this conclusion both to clarify the domains of study and scientific aims of cognitive science, psychology, and artificial intelligence, and to suggest some methodological principles for constructing intelligent machines.

## Acknowledgements

I owe much to JOSEPH SCHATZ for advice on this paper. I also benefit from discussions with NED BLOCK, JOHAN DE KLEER, MERRICK FURST, ALLEN NEWELL, DANA SCOTT, and RICHMOND THOMASON.

# Analysis of Contents

Section	Page
1. Recent work in artificial intelligence has shown a tendency toward rigor of proof and sharp definitions of concepts.	1
2. This critical examination must ultimately extend to the concept of mind itself. The aims of psychology.	1
3. Philosophical motives for such an inquiry: questions of mind and matter, free will and determinism, and effective computability.	2
4. Task of the present work.	2

## I. Views of certain writers on the nature of mind

### *Are there thoughts?*

5. Immaterial minds and material bodies. DESCARTES.	3
6. Material minds and bodies. EPICTETUS, HOBBS, WATSON. Psychological Behaviorism.	3
7. Philosophical Behaviorism. SKINNER. RYLE.	4
8. A fatal flaw in philosophical behaviorism.	4
9. The retreat from philosophical behaviorism. Identity theories.	4
10. Flaws in identity theories. PUTNAM. DENNETT. FODOR.	4
11. The retreat from identity theories. Functionalism.	5

### *Have thoughts content?*

12. Theories of the functions of thoughts. Naive and rigorous psychology.	6
13. Wide and narrow psychological states. Methodological solipsism.	6
14. The role of content in psychological theories. Comparisons of content. Propositional attitudes.	7
15. Thoughts as attitudes towards or expressing propositions. BRETANO and FREGE. On graspability of propositions. PUTNAM. KAPLAN.	7
16. Thoughts as attitudes towards sentences or representations. QUINE. FODOR. Universal, individual, eternal, and evolving languages of thought.	8
17. Functionalism emaciated by psychological solipsism.	9
18. Possible-world interpretations of thought. Anti-reductionism and the return to functionalism.	9

## II. Views on the nature of machines

### *Models of computation*

19.	Problem decomposability, reliability of mechanisms, and the historical simplicity of machines.	11
20.	The modern theory of machines. Computers. Effectiveness and Narrowness.	11
21.	Finite state machines. Monolithic machine states.	11
22.	Turing machines. Monolithic controller but fragmented tape.	12
23.	Random access machines. Decomposability assumptions, subroutines, and linking loaders.	12
24.	Functional programming languages. Full decomposability but no structure sharing.	12

### *Functional specifications in programming*

25.	VON NEUMANN'S architecture for modern computers and programming methodology.	13
26.	Program specifications as functional and ecological specifications.	13
27.	Syntactic and semantic verification of machine specifications. Relating specifications of parts to specifications of the whole.	14

### *Functional specifications in artificial intelligence*

28.	The differing roles of functional specifications in artificial intelligence and in ordinary programming. Programs as research vehicles vs. programs as results.	15
29.	The ordinary audience for specifications is the human programmer. Real-world interpretations of machines.	15
30.	The artificial intelligence audience for specifications is the machine itself. Adaptive-ness, self-programming, and routine updating of the state of mind. Typical forms of self-specifications.	15

## III. The concept of mind

### *Convergence of the theories*

31.	Convergence of narrow psychological theories in the philosophy of mind and narrow self-specifications in the practice of artificial intelligence.	17
-----	---	----

*Possible minds defined*

32.	The set of possible minds is the set of narrowly realized theories.	17
33.	Advantages of the definition. Neutrality of cognitive science on matters of psychological ontology, complexity, effectiveness, and determinateness.	17
34.	Recognition of the generality of the definition.	18

IV. Conclusion

35.	Methodological consequences of the definition. Expository advantages of the viewpoint for artificial intelligence.	20
36.	Advantages of the viewpoint for constructing intelligent machines. Constraint-based programming systems. Thinking as a process of narrow self-specification and self-interpretation.	20
References		22

*Archaischer Torso Apollos*

by R. M. Rilke  
translated by R. Bly

Wir kannten nicht sein unerhörtes Haupt,  
darin die Augenäpfel reiften. Aber  
sein Torso glüht noch wie ein Kandelaber,  
in dem sein Schauen, nur zurückgeschraubt,

sich hält und glänzt. Sonst könnte nicht der Bug  
der Brust dich blenden, und im leisen Drehen  
der Lenden Könnte nicht ein Lächeln gehen  
zu jener Mitte, die die Zeugung trug.

Sonst stünde dieser Stein entstellt und kurz  
unter der Schultern durchsichtigem Sturz  
und flimmerte nicht so wie Raubtierfelle;

und bräche nicht aus allen seinen Rändern  
aus wie ein Stern: denn da ist keine Stelle,  
die dich nicht sieht. Du musst dein Leben ändern.

We have no idea what his fantastic head  
was like, where the eyeballs were slowly swelling. But  
his body now is glowing like a gas lamp,  
whose inner eyes, only turned down a little,

hold their flame, shine. If there weren't light, the curve  
of the breast wouldn't blind you, and in the swerve  
of the thighs a smile wouldn't keep on going  
toward the place where the seeds are.

If there weren't light, this stone would look cut off  
where it drops clearly from the shoulders,  
its skin wouldn't gleam like the fur of a wild animal,

and the body wouldn't send out light from every edge  
as a star does ... for there is no place at all  
that isn't looking at you. You must change your life.



*for Gerald Jay Sussman,  
whose life is contagious*

## Introduction

If today one asks "What is the brain?" one receives volumes of material from the neuroscientists. This material may not be very complete, nor terribly revealing about the operation of the brain, but it is a beginning one can hope to continue. On the other hand, if one asks "What is the mind?" or "What is thinking?" one receives today little information more certain than that supplied by the quaint theories of antiquity. On these latter questions one finds not just volumes of material, but volumes of material for each of many diverse theories of mind. Philosophers alone supply dozens, but they are not exceptional, for so do the psychologists, linguists, decision-theorists, artificial intelligence researchers, and novelists. These and other answers are in turn assumed by the social sciences, each of which depends on a conception of man for its formulation, and by moral thinkers in philosophy, theology, and politics for similar reasons. What's a body to think?

Is it not a scandal that such common notions as mind and thought enter so prominently in our understanding of the world, yet find so little definiteness in the sciences they underlie? How are the sciences to look for facts, to measure, analyze, or construct minds without a clear notion of the object of their study? Without a clear conception of mind, how can the sciences tell if they are studying minds or if they have accidentally drifted to studying something else instead? How, indeed!

In spite of the great wealth of theories extant, my purpose here is to present yet another as an introduction to a number of works in preparation on artificial intelligence. The concept of mind is ultimately a philosophical or metaphysical subject, and as a consequence my inquiry begins in largely philosophical terms. I approach the subject from backgrounds in mathematics and artificial intelligence, and use of mathematical, logical, and computational methods soon join philosophical methods in my investigations. I believe the use of both philosophical and mathematical tools necessary to this project, even if some in artificial intelligence and other fields find these tools distasteful. If successful, these results will inform researchers in artificial intelligence of the logical foundations and implications of their techniques, and will inform philosophical and mathematical logicians of both psychological applications of their theories and possible areas for new mathematical developments. I permit myself the hope that even the philosophers, if they examine what I have written without prejudice, will find in it something of use to them.

The best of tools are powerless, however, if the nature of the investigation is not clear. Sadly, there is some confusion about the aims of artificial intelligence and the other cognitive sciences. One finds views of artificial intelligence ranging from "frontier applications of computers," to "knowledge engineering," to "making machines act more intelligently," to "making machines simulate human thought," to "making intelligent machines." A glance at introductions to cognitive science reveals scant definition beyond a melding of artificial intelligence and psychology, together with hopes that maybe bits of philosophy, linguistics, education, neurosciences, etc. will make their way into the field. This is a recipe for a stew, not a statement of the scientific aims of a field. Artificial intelligence and cognitive science may in fact be or turn out to be all these things, but these must be consequences of the pursuit of yet-unarticulated scientific aims rather than the aims themselves.

Clarity about scientific aims is especially important for the present work, for our investigations will not make much sense or will be misconstrued without a proper understanding of their intended contribution. The title of this work is *The Foundations of Psychology*, yet one will find virtually nothing from modern psychological theories here. This paper is not so much about the modern discipline of Psychology as it is about all possible organizations for minds, i.e. psychologies. These are the objects of study of cognitive science, as I see things. The disciplines of psychology and artificial intelligence are simply particular sub-disciplines of cognitive science with special interests to pursue. I share this last conclusion with the views mentioned above, but the former conclusion, upon which to me the latter is based, seems much less widely held, if others hold it at all. To avoid misconstrual of these investigations, the remainder of this introduction sets out the assumed scientific aims of cognitive science, psychology,

and artificial intelligence.

Cognitive science is the study of all possible minds. It may not be possible to set out in advance a definite class containing all possible minds, just as biologists have had to abandon all definitions of "living things" and adopt an accommodating approach to newly discovered life forms. Nevertheless, this paper formulates a definition of what minds are as an initial foundation for their study. To jump ahead of ourselves, we view minds as narrowly realized theories, so that minds are not natural objects but theory-relative instead. The theories realized as minds are psychologies (not to be confused with the discipline Psychology), ways of viewing the organizations of minds. This definition can hardly mean much now, but we do assume some definite range of objects as minds. The task of cognitive science is to discover classifications of minds so that each mind can be uniquely characterized in terms of the system of classifications, and so that identically characterized minds are isomorphic in some natural sense. In other words, the aim of cognitive science is to characterize the equivalence classes of possible minds. Some classifications of minds will involve the sorts of constituents (e.g. beliefs, pains, etc.) from which mental states are constructed. In addition to internal structural classifications, other classifications involve relations between psychologies and other things. In the case of relating one psychology to another, one has classifications involving homomorphism, embeddability, or compatibility, of one psychology being a form of another. For example, if humans typically realize a general psychology  $\psi$ , then a psychology  $\psi'$  realized by a particular normal human, will be a more detailed version of the general psychology, and so will admit a homomorphism onto  $\psi$ . In the case of relating psychologies to non-psychological entities, the principal question is whether a particular psychology  $\psi$  is realizable in entities of class  $E$ . For example, the task of the familiar discipline of Psychology is to find psychologies that can be and typically are realized in human beings. The task of artificial intelligence is to discover "interesting" (e.g. human-like) psychologies which can be realized in Turing-equivalent machines.

Human beings and Turing-equivalent machines need not exhaust the range of entities in which to realize psychologies. Many interesting psychologies may lie beyond the realm of what is realizable given the physics of our universe. Consideration of this possibility may seem strange to those steeped in traditions involving empirical psychology and CHURCH'S thesis, but limiting the scope of cognitive science by the laws of physics is mistaken. The questions of cognitive science have content independent of the particular characteristics of our universe. This can be understood in several steps as follows.

Individual psychologies are particular theories of mind, precise specifications of mental organizations. If a psychology is a formal theory of a mental organization, realizations of this psychology are models of the theory. There is no normative or descriptive content in the notion of psychology itself, just the notion of a theory and its possible models. However, we may view psychologies from both normative and descriptive perspectives. We can ask for the psychologies describing some individual human, for example, and we can say some psychology specifies normal human mental organization. When viewed normatively, psychologies are competence theories in CHOMSKY'S sense, where realizations of the normative psychology are mentally "competent" agents, and realizations of differing psychologies are mentally "incompetent" agents (although their incompetence may involve being stronger as well as weaker in their faculties, unavoidable supercompetence considered a form of incompetence at meeting normative limitations).

CHOMSKY'S idea of competence theories has influenced much work in the cognitive sciences, but the applications of this idea may have been unnecessarily limited in comparison with the more basic notion of normative theory due to the context in which the idea was introduced. In the beginning of *Aspects of Syntax*,<sup>1</sup> CHOMSKY develops his competence-performance distinction as a tool in explaining how a finite mind can use an apparently infinite language, to say how a grammar makes possible the "infinite use of finite means." CHOMSKY motivates the notion of competence with a picture of an idealized speaker as one free of all memory limitations and free of certain sorts of computational limitations such as distractions, shifts of attention and interest, and random or characteristic errors.

---

<sup>1</sup>[CHOMSKY 1965]

Since CHOMSKY speaks from a linguistic tradition which gathers yes-or-no judgments of grammaticality from subject speakers, it is easy to assume he assumes recursiveness of spoken languages. But to extend the requirement of recursiveness to the idea of competence seems unwarranted.<sup>2</sup> "Finite means" might be finitely axiomatized second-order theories like arithmetic, beyond the pale of recursive enumerability. Once one widens one's interest from human languages to cognitive science, limitations of recursiveness can be viewed as merely a sort of characteristic error for the speakers of a non-recursive language. The aim of a normative or competence theory is to give an ideal against which to measure the performance of supposed practitioners. There is nothing in the notion of ideal *simpliciter* which entails recursiveness or even recursive enumerability. We should be able to consider an ideal consisting of just the true sentences of arithmetic. Of course, we can prove that no Turing-equivalent speaker can achieve this ideal of competence, but that is nothing new in the world. Suppose, for example, we try to formulate weight-lifting competence and performance. It seems plausible to take the notion of weight-lifting competence to be that one can lift barbells of any weight. Of course, there are many "processing" limitations idiosyncratic to all humans. Moreover, one can prove from physical laws that there are weights humans can never lift. Is this provable physical incompetence somehow specially different from provable mental incompetence? It seems unlikely that CHURCH'S thesis is anything but an empirical fact, just like the ordinary laws of physics. GANDY, for instance, has outlined a proof of CHURCH'S thesis from physical laws.<sup>3</sup> The apparent equivalence of effective computability and recursive enumerability may be nothing more than an amazing coincidence of our universe. But if the limitations of CHURCH'S thesis are empirical rather than logically necessary, there might exist universes in which machines have super-Turing powers of computation, and excluding these universes and their inhabitants from the domain of normative theories is nothing but chauvinism.

By virtue of their use as normative theories, we must conclude that psychologies need not be actually realizable to be of interest to cognitive science. This range of concern frees cognitive science to calmly study controversial assumptions. Cognitive science studies in part the range of psychologies realizable by machines limited to the effectively computable, where what may be computable effectively by machines may vary with the laws of physics. Artificial intelligence will be the main research vehicle for studying mechanical intelligence in our universe. Cognitive science also studies the full range of psychologies realizable by arbitrary physical systems. In some universes, these may include non-mechanical psychologies, and in others, these may be the same as the mechanical psychologies. For example, human psychologies may or may not be mechanically realizable in our universe, but this is no cause for heated debates to impede the progress of cognitive science. Mechanists will relate their studies to the comparable but independent questions of artificial intelligence, and non-mechanists will not be so bothered unless it is to demonstrate the non-effectiveness of some aspects of human psychologies.

In addition to this argument from anti-chauvinism, one might arrive at the proposed scope of cognitive science as one more fruitful scientifically than one limited to our universe. Arbitrarily limiting the range of acceptable psychologies may unnecessarily prevent discovery of important facts about uncontroversial psychologies. The more generally posed problem may admit an easy solution even when the special case presents intractable difficulties. For example, real and complex numbers may be "unreal" in some sense in which integers are not, but some facts about integers and integer functions (e.g. the prime number theorem) are obtained most easily as facts about real and complex numbers and functions, and can be obtained only with extreme difficulty as facts about integers alone. In the same way we may most easily see things about the minds of men by looking, so to speak, to the mind of God. Limiting the domain of cognitive science to the boundaries of our own universe and excluding the trans-computable is analogous to limiting mathematics to the integers and excluding the transfinite. This analogy has a moral. CANTOR led mathematicians into a paradise many will not abandon due to the rich harvest they find there. Other mathematicians stay outside to see what fruit hang on branches

---

<sup>2</sup>[THOMASON 1979] discusses some problems with the notion of competence theories for cognitive scientists wedded to the assumption of recursiveness.

<sup>3</sup>[GANDY 1980]

reaching over the walls. In comparison, cognitive scientists have hardly begun to taste the fruit of their field. Perhaps our paradise grows only sour grapes, but we must taste them to see.

§1. After proceeding for centuries without exceptional standards of rigor, philosophy of mind within this century has been the beneficiary of three great boons. The first of these was the development, beginning primarily with FREGE, RUSSELL, and HILBERT, of modern mathematical logic and metamathematics.<sup>4</sup> The second boon was the development, beginning with WUNDT, FREUD, and WATSON, of modern scientific psychology.<sup>5</sup> The third boon was the development, with the advice of TURING and VON NEUMANN, of the general-purpose electronic computer.<sup>6</sup> Mathematical logic has provided the language for precise formulations of problems of mind and their solutions. Scientific psychology has provided many careful observations of human behavior and performance to account for in theories. And the high-speed electronic computer has provided the means for experimentally investigating the consequences of partial and comprehensive theories of mind in ways that were previously infeasible. Together these three developments have stimulated the new field of artificial intelligence. With the intellectual and practical tools at hand, artificial intelligence promises the most precisely formulated and visibly detailed theories of mind yet developed. Where previous centuries served as grounds for battles between imprecise, ill-understood proposals, battles allowing little hope for eventual comprehension of or agreement on the nature of the issues at stake, the new tools permit far greater clarity in formulating the concepts and structures of theories, so that parties to debates at least can agree on the issues, if not on the answers. Properly applied, the new tools ruthlessly expose inadequacies of theories previously hidden behind the general vagueness of formulations of the theories, and this shows up in the aphorism of artificial intelligence which states that about the nature of mind, the first twenty obvious ideas are wrong. More clearly than ever before, the new tools prove that in psychology simple moral conviction, even supported by a mass of successful applications, is not enough. Rigor in formulation must accompany conviction if the field is to visibly progress.

§2. Precision in formulating psychological theories leads quickly to a need for clear conceptions of the nature of mind and the nature of thinking. In the past, psychological theories have all too often lacked clarity through vagueness about the object of the theory. Routinely, investigators would propose hypotheses in great detail about some particular component of mind or thought, while merely alluding to an overall structure never made precise. This course holds two dangers: first, that not being closely tied to a definite conception of mind, the component theory will wander off to become an abstract plaything of its creator; and second, that the component theory will become vacuous by pushing all important theoretical burdens onto the vaguely defined theory of the whole. To avoid these dangers, we must strive for as much precision in setting out concepts of mind as in proposing theories of particular mental components.

If first aim of psychology must be to set out at least skeletal theories of mind, lest any empirical or detailed studies of supposed mental components risk triviality or irrelevance, then the first step toward this aim must be to understand the possible forms of psychological theories, to understand the range of

---

<sup>4</sup>See, for example, [FREGE 1884], [WHITEHEAD AND RUSSELL 1910], and [HILBERT 1900].

<sup>5</sup>See, for example, [WUNDT 1874], [FREUD 1895], and [WATSON 1914].

<sup>6</sup>See [GOLDSTEIN 1972] and [RANDELL 1975] for histories.

possible conceptions of mind. Care is required in this enterprise, for historically there has been frequent temptation to confuse the theoretical tools used to formulate theories with the theories themselves, and this confuses the aims of psychology with the aims of mathematics, logic, computer science, and other fields. For example, logic is not part of psychology, and neither is psychology a part of logic. Nevertheless, the temptation is frequently great to transfer the methods of logic to the methods of thinking and *vice versa*. Just as Modus Ponens is well established in logic, many non-deductive inferential processes are well established in psychology, and confusing the logical notions of entailment and proof with the psychological notions of argument and thinking is just wrong. Yet numerous students of mind make this confusion, and in consequence many works in artificial intelligence place psychological burdens on logic, faulting it for remaining silent on properly psychological problems, and mistakenly turning away from the clarity of formulation logical tools provide. Complementing this, other works view the purpose of thinking as the production of new knowledge from old, to which logical deduction is well suited if anything is. This transference of purpose oversteps the proper role of logical rigor in psychology. The purpose of thinking is insight, not knowledge. The aim of reasoning is not merely to ensure the truth or drawing of a conclusion, but more fundamentally, to afford insight into the dependence of conclusions upon one another. From certain sorts of interdependencies of conclusions, we can tell that one conclusion will be true if others are; and from other sorts of interdependencies, that one conclusion will be held by the agent if others are. On the face of it, these are different sorts of relations between conclusions. It is not difficult to think of further sorts of interesting relations between conclusions, such as confirmation or falsification in inductive reasoning. In fact, each of the disciplines brought to bear in formulating a conception of mind may have its own special sorts of interdependencies, but these must not be confused with each other. The separate disciplines must enhance each other, not replace each other.

§3. Although I am led to ask these questions in the course of formulating psychological theories, philosophical motives too have prompted me to enquiries of this kind. Our understanding of the questions traditionally raised about the relation of mind and matter, of free will and determinism, and of the limitations of effective computability all must increase from consideration of the nature of psychological theories. We may, if lucky, arrive at answers to these questions, or at formulations of possible positions if not answers. Even if no conclusions result, at least the outlines of the mysteries will be clarified.

§4. Surprisingly, starting from these theoretical and philosophical questions about the nature of mind, we are led to formulate the same demand as which had arisen independently in the practice of artificial intelligence, namely that psychological theories be cast as sets of narrowly interpreted self-specifications; for only if a theory is suitably narrow in the range of its references can it be realized in humans or machines; and only if a theory can be used in isolation to reconstruct states of mind after changes does it facilitate the design of artificial agents. Our story is largely a tale of the two theories leading to this remarkable convergence.

## I. Views of certain writers on the nature of mind

### *Are there thoughts?*

§5. During most of the history of philosophy, men have been assumed to have thoughts. More often than not, thoughts and their vehicle, the mind, were taken to be different somehow than the ordinary stuff from which our bodies are composed. A high point in this train of thought was the skeptical argument of DESCARTES.<sup>7</sup> By elementary considerations, DESCARTES managed to convince himself that his mind existed. He required substantial additional hypothesis and argument to convince himself that his body and other ordinary things existed as well. Based on these considerations, DESCARTES divided the world into two parts: the realm of matter, from which our bodies are composed, and the realm of an immaterial sort of substance, from which our minds are composed. For a long time after DESCARTES, the principal themes in discussions of mind consisted of speculative theories about the relation and possible interaction of immaterial mind and material body. These theories were hampered by the growth of the physical sciences, for as more and more of the universe came under the domain claimed by physics, it became increasingly difficult to supply a plausible account of how an immaterial object, not subject to physical law, could influence or be influenced by a material object subject to those laws.

§6. Fortunately for the growth of a somewhat less speculative psychology, another train of thought was developing in the wings. Long ago, the Epicurians, taken with DEMOCRITUS'S ideas about atoms, proclaimed that all the world was composed of atoms.<sup>8</sup> Men's minds, no less than their bodies, were made up of atoms, and the behavior of minds could be explained like the behavior of bodies, in terms of the mechanical laws governing the motions of atoms. However, the Epicurians proposed no details of these explanations and laws. Instead, they simply claimed the soul was composed of atoms lighter and more mobile than ordinary atoms, so presumably accounting for the rapidity with which thoughts can progress while leaving no outward sign of change visible to observers.

After some delay, the Epicurean programme was taken up by HOBBS, who tried to give some explanation of the workings of the material mind. HOBBS took his models of mechanics from the new physics, which clarified the notions of forces and inertia. With these, he explained thought and imagination as the decaying motion of the atoms of the sense organs as they stimulate other atoms and are stimulated by new impressions. In spite of the possibilities of this programme, HOBBS was limited by contemporary ignorance about neurophysiology and non-naive psychological phenomena, so his theory remained largely speculative and general about mental mechanics.

After yet another delay, these ideas received their first detailed treatment. An experimental, scientific psychology developed under the impetus of WUNDT and others, and much information became available about psychological phenomena.<sup>9</sup> However, much of this information still depended on uncritical speculation in the form of introspective evidence. WATSON urged that psychology secure its foundations and either shore up introspection or avoid it altogether in favor of theories of overt and neurophysical behavior.<sup>10</sup> He took the latter course, and his programme, psychological behaviorism, be-

---

<sup>7</sup>[DESCARTES 1637]

<sup>8</sup>Much of my knowledge of these fragments of history comes from [BORING 1950] and [PETERS AND MACE 1967].

<sup>9</sup>[WUNDT 1874]

<sup>10</sup>[WATSON 1914]



came widely influential, and led to much progress in matters such as neurophysiology, neurochemistry, and neuropsychology.

§7. Given this development of ideas, it seems like a small step from psychological behaviorism to philosophical behaviorism. Where psychological behaviorism, in modest formulations at least, seeks simply to determine the laws of mind and thought and their realization in the body, philosophical behaviorism, as championed by RYLE and SKINNER, claims reducibility to overt behavior as well.<sup>11</sup> For the philosophical behaviorist, mental phenomena can not only be realized in terms of neurological hardware, but can be reduced to those functionings of that hardware, in the sense that all theories about belief, desire, inference, and action can be formulated purely in terms of the overt behavior of the body. In this view, mental entities and processes are simply unreal, and their use is as unscientific as phlogiston descriptions of combustion. RYLE, for example, sought descriptions of human behavior strictly in terms of overt acts, dispositions to behave, and changes in dispositions to behave, where dispositions to behave were formulated purely as physical states and physical laws.

Philosophical behaviorism eschewed the mind. Indeed, some of its adherents went so far as to deny any introspective self-awareness in their own cases: not a tactic unknown to philosophy, but one never before practiced on such a grand scale. Cartesian skepticism seems meagre compared to that of the philosophical behaviorists.

§8. Philosophical behaviorism is implausible, and for very simple reasons. Suppose one attempts to analyze an ascription of belief, e.g., "Fred believes his computer program has a bug" in terms of behavior and dispositions to behave. We can of course guess at predictive generalizations about Fred's behavior given this belief, such as "Fred is disposed to log in and debug his program," but for any of these predictions to have any plausibility, they cannot be formulated simply in terms of behavior and states of the world; they must refer to other of Fred's beliefs and desires that might influence his dispositions, e.g., "Fred believes the computer is down" or "Fred wants to take his time so as to increase his wages." Because predictions of Fred's behavior must refer to his thoughts about things as well as to the actual state of things, philosophical behaviorism is untenable.

§9. With Cartesian dualism indefensible, but with its prime alternative inadequate, philosophers developed several replacement theories which identified mental states and processes with the physical states and processes of the brain and connecting parts of the body.<sup>12</sup> These identity theories, as they are known, hoped to analyze beliefs and desires, for example, as physical predicates of brains, thus allowing the objections to philosophical behaviorism to be overcome by phrasing references to beliefs and desires as references to certain sorts of brain states. The identity theorists returned to a position much like that of HOBBS, although unlike HOBBS, they had a stronger physics and neurophysiology to draw upon for formulation and examples.

§10. Although the identity theories avoid the obvious problems with philosophical behaviorism, they suffer from inadequacies of their own. Curiously, these inadequacies are made all the sharper by the

---

<sup>11</sup>[RYLE 1949], [SKINNER 1957]

<sup>12</sup>See, for instance, [PLACE 1956], [FEIGL 1958], and [SMART 1959].

development of somewhat intelligent machines in artificial intelligence, although the inadequacies can be brought out without using those machines as examples.

Identity theories claim that all statements about mental phenomena are reducible to statements purely about physical brain states, with the implication that perhaps we are better off studying the brain states in psychology and leaving the mind alone. The first objection to this idea, at least to its implication, is that even if such a reduction is possible, it is useless in practice and so even theoretically uninteresting (at least to those who include practical power as a measure of a good theory). PUTNAM gives as example the analogous case of a description of a hole in a board in terms of the location and momentum of elementary particles.<sup>13</sup> If the hole is round, we cannot insert a square peg of equal cross-section area, and presumably this fact can be explained in terms of the locations of and forces between the particles comprising the board and peg. But such an explanation must be astronomically long, and must involve many details that are in some sense irrelevant to the important facts, namely, the roundness of the hole and squareness of the peg. DENNETT gives as another example chess playing computers.<sup>14</sup> Unless one is a chess master, one's best tactic in playing such a machine is to treat it as a rational agent with beliefs, desires, and knowledge of the game, for the actual sequence of computational steps the machine employs to develop its moves is incomprehensible (at least in any reasonable amount of time for individual humans). DENNETT describes this observation as the notion of "intentional stance": that no matter what other theories are possible concerning the realization of mental phenomena, we still must use a psychology formulated in mental terms, because mental ascriptions are the ones of practical predictive power (for humans, at any rate).

The identity theories also fall prey to a stronger criticism: that their supposed reduction of psychological physical laws cannot be done, even in principle, without severely gutting the notion of law. Suppose one has a psychological theory formalized in mental terms, a physical theory formalized in terms of particles, and fields, and, as the identity theorists suggest, a set of "bridging laws" or "bridging definitions" connecting the two theories. FODOR points out that because there are presumably many, many ways of realizing minds in matter (just as many physical objects can be used as money), the bridging statements cannot be laws unless they are so wildly disjunctive as to ridicule our ordinary conception of what it means to be a law.<sup>15</sup> FODOR continues by pointing out that if the bridging statements are reasonably informative or concise definitions, then the psychological theory that results cannot aspire to the title of law either — its applicability need not be wide, nor need its conditionals have counterfactual force.

§11. These criticisms of identity theories contain the seeds of a theory to replace the identity theories, functionalism. Functionalists contend that mental states are functional states; that is, what makes something a belief or a desire of the agent is the fact that that something plays a certain role in the processes of the agent, independent of any other properties possessed by its particular realizations in brains, computers, or ghosts.<sup>16</sup> Expressions of functionalist psychological theories take the form of physical theory plus a collection of Ramsey sentences: existentially quantified statements asserting the existence of something that bears certain relations to other things which, in light of the criticism of psychological behaviorism, will include some of the other entities asserted to exist by the Ramsey sentences. For example, a ten cent soda machine accepting either nickels or dimes might be described by the two sentences "There exists a state  $S_0$  such that receiving a nickle in  $S_0$  results in a change to state  $S_5$ , and receiving a dime in  $S_0$  results in dispensing a soda and remaining in  $S_0$ " and "There exists a state  $S_5$  such that receiving a nickle in  $S_5$  results in dispensing a soda and a change to state

---

<sup>13</sup>[PUTNAM 1975C]

<sup>14</sup>[DENNETT 1978]

<sup>15</sup>[FODOR 1975]

<sup>16</sup>See, for instance, [PUTNAM 1975B] and [FODOR 1968].

$S_0$ , and receiving a dime in  $S_5$  results in dispensing a soda, dispensing a nickle, and a change to state  $S_0$ ." Functionalism thus permits some separation between the development of a psychological theory and the delineation of the ways in which the theory may be realized in humans or machines. However, functionalism is not without problems of its own, and some of these concerning surprises about what sorts of things can and cannot be functional states will arise in the following.

### *Have thoughts content?*

§12. Even if one accepts functionalism as a hopeful theory of the mental, there are many serious difficulties that must be overcome. The first of these is the nature of thoughts as conceived in functionalist terms. To understand this problem, we must examine the functions thoughts play in mental activity.

To begin a functional characterization of the role of particular thoughts in thinking, we can take as a crude approximation all the truisms about rational thought and action. For example, a functionalist psychology might have specifications like the following: (1) "If agent  $A$  desires  $p$  and believes  $M$  a method for achieving  $p$ , then  $A$  desires to do  $M$ , other things being equal," and (2) "If  $A$  desires to do  $M$ , and believes doing  $M$  is better than doing any other action  $A$  desires, then  $A$  does  $M$ ." Our lazy programmer Fred might desire to earn a living without the necessity of toil, might believe that being slow to debug his programs a way of doing this, and might in fact be slow in his chores thinking this the easiest way of life. Of course, specifications (1) and (2) are woefully simplistic, and any careful treatment would have to be based on a better theory of rationality and practical reasoning, but the example serves to illustrate the basic idea.

§13. In addition to basing functional specifications for the mind on sound theories of thinking and doing, we must take care to ensure that the specifications are suitably narrow in scope. For example, we might consider a specification like (1') "If  $A$  desires  $p$ , and  $M$  is a method for achieving  $p$ , then  $A$  desires to do  $M$ , other things being equal." The problem with (1') is that it puts a measure of omniscience into the specifications, omniscience that offends our notion of what a psychology should involve. While one might give such specifications in ecological theories of how well agents fit into their environments, for psychology proper we eschew any specifications which do not express what PUTNAM terms "psychology in the narrow sense."<sup>17</sup> Psychology in the narrow sense is concerned only with narrowly realized psychological theories, where narrowly realized psychological theories refer only to the agent's personal mental structures, and do not refer to any correspondence or lack of correspondence between the agent's thoughts and its environment, such as the truth of its beliefs or the physical possibility of its desires. Narrowness is a property of realizations or interpretations of psychological theories, rather than a property of the theories themselves, since a particular theory might admit interpretations stepping outside the agent as well as strictly internal interpretations. For example, to return to the ten cent soda machine introduced above, the functional specifications given there were not narrowly realized. Those specifications referred to dispensing of sodas, objects presumably beyond the machine's ken. Ordinary soda machines (at least the ones I have lost money to) only send signals to the dispenser for its jaws to open, so allowing a soda to fall if one is there, but failing to dispense a soda (unbeknownst to the machine) if the soda rack is empty or jammed. To describe the machine's structure narrowly, we must replace references to sodas by references to signals sent from the cashier to the dispenser. FODOR terms such restrictions of narrowness on acceptable interpretations of psychological specifications "methodological solipsism."<sup>18</sup>

---

<sup>17</sup>[PUTNAM 1975A]

<sup>18</sup>[FODOR 1981]

§14. Adherence to the methodology of psychological solipsism requires that we re-examine our naive psychological theories to excise all externalities. The example specifications (1) and (2) given above seem to refer to the content or meaning of the agent's beliefs, desires, etc., but notions of meaning or interpretation are usually external relations of the agent to the world, hence the use of such notions in our psychological theory is suspect. We could immediately excise these notions from our psychology, but this would be hasty since there are many reasons why we would like to have some notion of content available for thoughts. The agent must be able to compare or distinguish its thoughts, one from another, lest its psychology be completely trivial and implausible. Some notion of content is also required for comparing thoughts of the agent at different times, and with the thoughts of other agents. We clearly can make some powers of discrimination between thoughts available to the agent simply by relying on a formal syntax of thoughts, but it is less clear whether these discriminations will match or can match the discriminations needed in diachronic and inter-agent comparisons. These latter comparisons are more clearly dependent on interpretations external to the agent, but all of these cases call for careful examination.

Before considering the answers that have been proposed for these questions, we introduce a new bit of terminology. Functional specifications of psychologies have thoughts playing certain roles in the agent. We can turn this relation around and for convenience say that the agent bears a certain relation to its thoughts. Indeed, this latter phrasing suggests some of the motivation for RUSSELL'S "propositional attitude" terminology, since the agent-thought relation can be viewed as a certain stance or attitude taken by the agent toward a possible thought. In this phraseology, the requirement of methodological solipsism becomes the requirement that thoughts are graspable, that is, some how manipulable and determinable by the agent itself. The agent presumably bears some relation to the real world, but if the agent is only a small part of the world and the weakest forms of skepticism are justified, the agent cannot fully grasp its relation to the external world.

§15. The first suggestion we consider about the nature of thoughts holds that the content of thoughts are abstract entities called *propositions*. The term "propositional attitude" stems from this orientation. Propositions as the content of thoughts comes as a suggestion from two quarters: BRETANO'S theory of intentionality and FREGE'S theory of language.

BRETANO distinguished mental states from physical states of an agent with the idea of intentionality.<sup>19</sup> He claimed that all acts of consciousness, such as beliefs and desires, are *directed towards* or *about* some object or objects. For example, a hatred of paperclips is an act of consciousness whose intentional objects are all paperclips. In the case of the belief or desire that some condition obtains, BRETANO maintained that the object of the belief or desire is a proposition.

From another orientation, FREGE distinguished the sense of names or terms from their reference.<sup>20</sup> To use the celebrated example, the terms "author of *Waverley*" and "Scott" have different senses, but refer to the same person. FREGE claimed that thoughts express propositions as their sense, and may be either true or false as their reference. (Actually, this takes liberties with FREGE'S conception, but the details are not crucial here. Moreover, he used a term translated as both "thought" and "proposition", calling what we call thoughts "ideas" or "concepts".)

These suggestions have attracted many critics and defenders, but the main point we note about these suggestions is that such propositions cannot be part of a suitably narrow psychological specification for agent's structure, i.e. propositions are not graspable. PUTNAM observed that one can distinguish between real and apparent propositional content in beliefs, just as earlier skeptical arguments distinguish

---

<sup>19</sup>[BRETANO 1874]

<sup>20</sup>[FREGE 1892]

between the real and apparent truth of beliefs.<sup>21</sup> PUTNAM'S examples are rather involved, but KAPLAN gives the simpler examples of thoughts involving indexicals, i.e., implicit self or temporal reference.<sup>22</sup> For example, on Monday I think the thought "Today is beautiful." In the propositional theory, I can think exactly the same thought on Tuesday by thinking "Yesterday was beautiful." But this identity of actual content between these thoughts is indiscernable to me, since I might lose track of the days and think "Yesterday was beautiful" on Wednesday, thus actually expressing a proposition about Tuesday while thinking I am merely reminiscing on Monday's glories. Thus what proposition I actually express with my thought is not determined by my narrow psychological state.

§16. One can back off from propositions with the suggestion that thoughts are the relation of the agent to sentences. This suggestion has been motivated both in terms of an aversion for introducing abstract objects like propositions, and in terms of seeking graspable thoughts. Rather than introduce propositions as the meaning of sentences, and having propositions be true or false depending on circumstances, QUINE recommends simply having the sentences themselves be true or false in circumstances directly, rather than indirectly through the medium of propositions.<sup>23</sup> And rather than have an agent's thoughts involve grasping ungraspable abstractions, FODOR recommends having thoughts involve realizable representations, such as electrical patterns in the brain, markings on paper, or data structures in computer storage.<sup>24</sup> While this suggestion avoids the pitfall of a non-narrow psychological specification for the structure of mind, it raises other serious questions.

The first problem for the sentential view is that of the language of thought (sometimes called "brain-writing" or "Mentalese"). Since the sentences involved in thoughts are not just abstract entities but members of some language or representational system, they have a concrete syntax, and this places the sentential theory on the horns of a dilemma. Either the language of thought is common to all agents at all times, in which case diachronic and inter-agent comparisons of thoughts can be made within the confines of solipsistic psychology, or else the language of thought can depend on the temporal development of individual agents, in which case functionalism's claims to general applicability are sabotaged.

The hypothesis of a universal language of thought is difficult to accept for many reasons. The first objection is that if there are several (actual or possible) species of agents, the hypotheses claims a universality that can only be defended on grounds of cognitive necessity, on grounds that the very nature of successful or rational thought and action entails, by means of physical or computational necessity, the features of Universal Mentalese. While it seems plausible (perhaps weakly so) that some general features of language must be forced by the necessities of the task, it requires considerable demonstration that all features of the languages must be the same. One might claim that all agents of a certain species begin with the same innate language, or grow to accept the same language through learning of cultural conventions, but these suggestions are suspect on the grounds of simple genetic and educational variability.

But if the hypothesis of a universal language of thought is implausible, the acceptance of individual evolving languages is not without theoretical difficulties or unpleasant consequences. Theoretically, one may have to retreat from the view that a language has a definite set of sentences as members, for actual questions of membership may be infeasible or may change the language, so challenging the sensibility of determining the language by combining hypothetical judgments. Practically, acceptance of this hypothesis means that the accuracy of memory becomes much more problematic, since sentences

---

<sup>21</sup>[PUTNAM 1975A]

<sup>22</sup>[KAPLAN TBP]

<sup>23</sup>[QUINE 1970]

<sup>24</sup>[FODOR 1975]

brain-written in the past may play a different role in the agent's language in the future, if they still are in the language at all. This entails either continuous rephrasing of memories when the language changes, or loss of the power to interpret memories. This last possibility is not without attraction, since it fits well with Piagetian-style theories of psychological development, in which the conceptualization of the world changes radically as a child grows. I can hardly recall a thing of my childhood: perhaps I still have all my brain-records, yet cannot recall or make sense of them any more in terms of my present mental language.

§17. Psychological solipsism seems to force us to accept individual, evolving languages of thought with their attendant temporal and contentual ambiguities for the agent. This conclusion seems in some ways to defeat the motivation for functional specifications of psychologies. For example, recall specification (1): "If agent *A* desires *p*, and believes *M* a method for achieving *p*, then *A* desires to do *M*." Since narrowly interpreted psychology cannot refer to actual content, we cannot include specifications like "*A* desires *p*", since that refers to our interpretation of *A*'s attitude. We cannot even say instead "*A* thinks that *A* desires *p*", since that merely shifts the problem to the accuracy of *A*'s introspective beliefs. Instead, we are driven to write instead things like (1''): "If agent *A* incorporates *A*-now-Mentalese sentence *S*<sub>1</sub> and *A*-now-Mentalese sentence *S*<sub>2</sub>, then *A* incorporates (or adopts) *A*-now-Mentalese sentence *S*<sub>3</sub>," where we must write the concrete *A*-now-Mentalese sentences *S*<sub>1</sub>, *S*<sub>2</sub>, and *S*<sub>3</sub> in our theory, since we cannot refer to their external interpretation. But does such a theory fit our intentions for a psychological theory? We started by thinking we could sharpen up our naive beliefs about psychology by expressing them as a functional theory, but we find that if we want a narrowly interpreted theory, we get a purely formal set of specifications like (1'') for each individual agent at each atomic interval of time. Such a theory may leave open some questions about the details by which the agent actually realizes all these sentences and computes with them: but it certainly seems to have abandoned the generality of realization motivating functionalism in the first place. It may be possible to salvage these ideas by formulating general psychological theories which describe both universals of a species's psychology and which postulate the existence, for each agent at each time, a personal language instantiating or approximating the general one, but we have none to propose at this point.

§18. Another suggestion about the nature of thoughts is less easy to place within a non-behavioristic theory of mind, and that is the idea of possible-world interpretations of propositional attitude ascriptions.<sup>25</sup> In this approach one says *A* believes *p* if *p* is true in every possible world compatible with what *A* believes. This may seem circular, but one can take as primitive *A*'s acceptance or grasping of a set of possible worlds, among which *A*'s actual world is supposed to lie. With such a primitive conception, talk about beliefs is reduced to talk that does not involve beliefs in an attractive way.

However, this approach suffers from several difficulties. First, it entails that an agent's beliefs are closed with respect to logical consequence. It is certainly not obvious that (or if) the functional specifications for beliefs of any sufficiently rigorous intuitive psychology entail deductive closure, yet such closure seems to be an inescapable consequence of the possible-world approach. A second difficulty is that the earlier arguments about the ungraspability of meaning would seem to have similar force for the graspability of a set of possible worlds. If these problems are not enough, the third difficulty with the possible world approach arises when we try to extend the idea to other sorts of attitudes. If we take exact analogies for interpreting ascriptions of desires, hopes, fears, angers, etc., we get the same closure problem. But it is a strange psychology which does not allow for conflicting yet limited set of desires (or beliefs, for that matter). I can desire to have my cake and eat it too and still not desire nuclear war. In addition, a psychology should be sensitive not just to logical incompatibility between desires but also to

---

<sup>25</sup>See, for instance, [HINTIKKA 1971].

compatibility with respect to the agent's beliefs, yet the straightforward possible-world approach cannot incorporate such a notion. Instead, the only way open for making sense of possible-world compatibility seems to be by using the psychological theories to determine possible worlds as compatible sets of mental attitudes. And this is simply an acknowledgement of the anti-reductionist arguments mentioned earlier in support of functionalism.

## II. Views on the nature of machines

### *Models of computation*

§19. For most of history, man has had only simple machines. Machines had at most a few parts, often fixed rather than moving, and the invention of a new useful machine was a significant event for civilization. Part of the reason for this contrast with current times, in which vast numbers of new machines are constructed each day, must be that machines were less reliable in the past, and so if a complex machine was constructed from simpler ones, its chance of working was small. The big advances were not simply the construction of new machines, but of ones simple enough to be reliable, rather than Rube Goldberg contraptions. Would any king today pine for a mechanical nightingale?

The obstacle of unreliability is serious. One of our most powerful intellectual tools for understanding the world and for taking action is *problem decomposition*. Problem decomposition is the technique of breaking one problem up into several hopefully simpler subproblems together with a way of combining the subproblem solutions into a composite solution to the original problem. When I make a cake, I do not throw all the ingredients into the oven and expect a cake to appear. Instead, I break the task into first making the batter, then baking it, then making the frosting, and finally assembling the finished cake. Moreover, this trick works for more serious problems than cooking. But for problem decomposition to work, the composition methods must be effective. When most machines were not terribly reliable or precise, new machines could not be designed by straightforward problem decomposition because the composite machine would be hopelessly unreliable, if it worked at all. Until science and technology had progressed to the point where most machines could be built to be reliable, invention of useful new machines remained largely a matter of inspiration, luck, or natural analogy.

§20. The modern theory of machines developed with the advent of reliable machines, and has focussed mostly on computational devices, since modern computers and the specialized machines realized in them via programs are the most complex machines ever known. Yet though modern machines are more reliable and can be combined in more complex constructions, theoretical models of machines do not always reflect this. Instead, the earliest developed models of machines offered little sense of machines decomposable into parts, and concentrated solely on a notion of the machine as a whole.<sup>26</sup>

Whatever their inadequacies regarding decomposability, machine models fulfill the subsequently formulated philosophical demands about narrowness. Normal uses of machine specifications involve narrow interpretations of the machine states. The idea of machine is closely connected with the idea of effective calculability, and non-narrow interpretations of machine specifications make the machine operations non-effective. In fact, effectiveness is a stronger requirement than simple narrowness, and later we see what sorts of things might fill this gap.

§21. One of the first models of machines was that of the *finite state machine*, introduced at least as early as 1936 by TURING.<sup>27</sup> A finite state machine is simply a transducer of input strings into output strings, with a finite amount of memory. Each finite state machine  $M$  is completely described by two

---

<sup>26</sup>For presentations of many sorts of machine models, consult [MINSKY 1967] and [AHO, HOPCROFT, AND ULLMAN 1974].

<sup>27</sup>[TURING 1936]



functions  $F_M$  and  $G_M$ , such that if  $M$  is in state  $S_j \in \{S_1, \dots, S_k\}$  and receives input  $I_j \in \{I_1, \dots, I_l\}$ , it emits the output  $F_M(S_j, I_j) \in \{O_1, \dots, O_m\}$  and moves to state  $G_M(S_j, I_j) \in \{S_1, \dots, S_k\}$ . The finite set of states of the machine limits its memory capacity, so that there are severe restrictions on what can be computed with a finite state machine. Nevertheless, finite state machines are theoretically interesting and frequently useful in practice. Unfortunately, this model of machines provides a relatively poor foundation for understanding machines by means of problem decomposition. Each of the states of a finite state machine is atomic, so there is no overt sense of a finite state machine having parts or being constructed out of submachines. Indeed, one of the heights of the theory of finite state machines is a characterization of when a finite state machine is equivalent to the sum, product, or concatenation of smaller finite state machines. We must look further to find models of machines congenial to problem decomposition.

§22. Another early machine model is the *Turing machine*.<sup>28</sup> Turing machines are simply finite state machines which can read and write symbols from a finite alphabet on an infinite tape. The tape symbols form the input alphabet of the finite-state controller, and combinations of symbols to write and tape motion signals form the output alphabet. Although Turing machines put some of their state onto the tape, where decompositions can be observed by using separate areas of the tape for separate sub-computations, Turing machines do not facilitate problem decomposition much more than do finite state machines, since the tape controllers are just finite state machines with the difficulties observed previously.

§23. Machine models closer to the structure of modern computers were formalized in the *random access machine* (RAM) and *random access stored program machine* (RASP).<sup>29</sup> A RAM has a read-only input tape, a write only output tape, an infinite array of memory cells, and a program. Each tape square or memory cell is either blank or may contain an integer. The machine's program is a finite list of numbered instructions from a certain fixed repertoire which can make arithmetic computations and comparisons on the memory cells as well as specifying the number of the next instruction to execute. RASP's are just like RAM's except that the program is stored in the memory cells, so that the machine can modify its own instructions by altering the contents of the memory cells storing the program.

RAM's and RASP's are much better than finite state machines or Turing machines at facilitating problem decomposition. Because both the program and memory are broken into discrete components, RAM's and RASP's computing a combination function can be built more or less by concatenating and renumbering the programs and by relocating the memory segments used by each sub-machine to disjoint (possibly interleaved) components of the combined memory. These are essentially the ideas of subroutines and linking loaders so important in modern programming systems. Such combination of machines is possible because each of the operations of the machine changes only a bounded component of the machine's state, i.e. a couple of cells and the program counter. This means that all operations ignore almost all of the machine state, so that separate sub-machines ignore each other when combined, except in their desired communication channels.

§24. Models of computation moved most recently to a position of complete decomposability in the functional programming languages.<sup>30</sup> We must live with an unfortunate coincidence of terminology

---

<sup>28</sup>[TURING 1936]

<sup>29</sup>See [SHEPHERDSON AND STURGIS 1963] and [ELGOT AND ROBINSON 1964].

<sup>30</sup>See, for example, [BACKUS 1978].

between "functional" specifications of roles and "functional" programming languages. The notion of roles played by some object need not be the same as the mathematical notion of function. We try to minimize confusion by always referring to the latter notion as "functional programming." In the functional programming model, submachines realize self-contained, arbitrary type functions, and composite machines are constructed by various sorts of functional composition and application from sub-machines. Since no machines share any structure, there can be no interference when they are combined. This greatly facilitates problem decomposition. In fact, the principal theoretical difficulties with functional programming languages involve how to re-introduce shared structures and side-effects in a useful way. Shared structures are often important for economy of storage usage, and for economy of effort in updating a database common to a number of separate processes. Unfortunately, functional programming languages have gone too far in seeking decomposability, since structure sharing is outside the domain of a pure functional programming language, so that a less extreme position must be found.

### *Functional specifications in programming*

§25. Most actual computers in use today are organized something like RASP's in the so-called *Von Neumann architecture*. Reflecting this common architecture, most programming languages are based on the idea of combining procedures with local instruction and data sets communicating through global or shared data sets, although there are a few languages like LISP and APL which come close to the functional language conception. One important consequence of the typical structure of computing machines and programming languages, if it is merely a consequence and not more deeply intertwined, is the phenomenon of procedural thinking among programmers. Most programmers find it easiest to compose programs by conceiving or imagining sequences of operations that wind up with the intended result. Only after they have composed the program in this way do they, if ever, reflect on what the program is computing to explain the functional relationships between the pieces of information it manipulates. This do first, reflect later phenomenon may be simply a consequence of the intellectual culture in which the programmers were raised, or it might stem from the human mind being a decomposable system in which most of one's mental state is automatically conserved from one moment to the next, facilitating envisioning of individual actions. In any event, programming methodologists have had to develop ways by which conscious functional decomposition can be facilitated in programming. The popular methodologies ("structured programming," "stepwise refinement," etc.) are suggestions for how to break problems up into explicit subproblems, how to combine the separate sub-solutions, and how to ensure or check the correctness of any non-interference assumptions made in the process of combination. Of course, these methodologies are still more hints about how to think than recipes for programming. The topic of how to decompose problems into subproblems is still more the domain of artists than of technicians, as only incomplete heuristics have been articulated, and the mechanization of these is still part of artificial intelligence. In contrast, much more is known about ways of checking the correctness of non-interference assumptions. There are two issues here: one of how to state what a machine is intended to compute, as opposed to how it is supposed to compute, and one of how to relate the intentions for sub-machines to intentions for the whole machine.

§26. Although each computer program is a precise set of instructions written in an interpreted formal language, most programming languages provide no formal means for restating what the program's instructions are intended to accomplish. Most programming languages provide only a commenting facility, with which the programmer can attach to lines of program text comments in English (or some other natural language) to indicate what the individual program instructions mean in the larger scheme of the program. Unfortunately, these comments do not have equal force of specification as the program text,

but instead are completely non-operative. In part this is due to the informal language of comments, to the lack of a formal language within the programming language for stating comments. The result is that the program forms a uniquely privileged specification of the machine realizing it, the only description that really matters to the operation of the machine. Some programming methodologists abhor the discretion this privileged status gives programmers in deciding whether or not to document their programs, but the number of machine descriptions does not seem to be the crucial issue. While undocumented programs are often odious, the underlying problem seems one of lack of force of specification (either to the computer or to programmers) rather than one of discretionary use.

To remedy these deficiencies of programming languages, computer scientists have developed numerous formal specification languages, logical metalanguages of programming languages with which one can state the intended effects of a program (its ecological specifications) and the intended roles of program components in program operation (its internal functional specifications). The ecological specifications for the program as a whole connect to the internal functional specifications by means of ecological descriptions of the elementary sorts of program instructions. For example, one common type of instruction operates on certain data-structures so as to mimic arithmetic operations on numbers. Since almost everyone thinks of these instructions as arithmetic operations, or uses them for other purposes by means of arithmetic encodings, the axiomatizations of these instructions usually state the effects of the instructions' execution in terms of arithmetic operations on numerically interpreted data-structures. Similarly if computers had instructions mimicking the operations of sewing machines on cloths, one might axiomatize these instruction in terms of sewing operations on cloth-interpreted data-structures. Other sorts of instructions call for further fabrication of interpreting axioms.

§27. Complementing the notion of formal specification of machines is the notion of verification of machine specifications, in which one checks that the specification of a composite machine follows from the specifications of the sub-machines and properties of the combination method. That is, one asks whether the ecological specifications of program operation follow from its functional specifications plus the ecological axiomatizations of instruction execution effects. Verifications can be approached syntactically or semantically, either by giving a logic of programs for formally deriving relations between specifications, or by giving a model which simultaneously satisfies all the specifications. While the syntactic approach is more immediately amenable to mechanization (and many attempts at mechanization populate the literature), the semantic approach is more fundamental, since a theory of models must underlie any logic of programs. Models of external domains like arithmetic, symbol strings, cloths, etc. present no difficulties peculiar to computer science, but many sorts of instructions operate on the internal components of machine states, affecting by their actions the meaning of the instructions themselves.<sup>31</sup> Models for these sorts of instructions are much more complex and much less familiar than everyday models involving ships and sealing wax. This unfamiliarity may be part of why programming seems so hard to teach; in any case, this complexity presented serious obstacles to development of satisfactory models for program specifications. However, much progress has been made, the most striking advance being the models for type-free functional programming languages developed by SCOTT, PLOTKIN, and others.<sup>32</sup> These developments make doubly appropriate the term "functional specification," since they allow functions to be elements of the domains of models as well as parts of the relational structures of the models.

---

<sup>31</sup>See [ELGOT AND ROBINSON 1964] for instructive models of RASP's in these terms.

<sup>32</sup>See [SCOTT 1973], [PLOTKIN 1972], [BARENDREGT 1981].

### *Functional specifications in artificial intelligence*

§28. While formal specifications fill the literature and textbooks of computer science, on first glance they seem almost totally absent from the artificial intelligence literature, their appearance there being restricted to attempts at mechanizing the problem decomposition and specification verification processes for "automatic" (i.e. machine performed) programming. In actuality, however, formal specifications do play a significant role in artificial intelligence, but a very different role from that played in computer science. This difference in role stems partly from the nature of artificial intelligence research, where problems of formulation play so great a role. If the primary purpose of writing a program is to increase one's understanding of the psychology realized by the program, verifying the specifications of the program becomes a pointless activity compared to reformulating the psychology and the implementation to accord with the insights gained in the experiment. Most programs written in artificial intelligence are not meant to be solutions, but are meant to be rungs on the ladder of understanding, rungs which allow one to progress, but rungs to be discarded after use. This is quite a contrast with the usual situation in computer science, where the focus is on better algorithms for tasks with stable formulations. But however important this difference in the purposes of programming in the two fields, it is dwarfed by the difference between the audiences of specifications in the two fields.

§29. In computer science, the study of machine specifications focusses on how a human programmer can think about programs and their use. This means that the principal contribution of formal specifications and semantics is a way of interpreting the structure of machines in terms of their real-world meaning. A hand calculator or slide rule is interesting only because we interpret the data-structures it manipulates as numbers, and because we interpret its operations on these data-structures as arithmetic operations. We do not care about changes in machine states, flip-flops, or bit patterns. We care about numbers and arithmetic. Similarly, we do not specify a bank's accounting machine in terms of patterns of bits and their manipulation, but in terms of customers, their deposits, their withdrawals, and their balances. The interpretations of interest are those relating the machine's state to the external world, and any discussion of the relations between machine states is merely part of a proof that the specifications of external interest are reflected in the structure and behavior of the machine. In our earlier language for describing psychological theories, the specifications of interest in most of computer programming are non-narrow, ecologically interpreted specifications which refer to circumstances external to the machine.

§30. In artificial intelligence, however, the study of machine specifications focusses on how the machine can think about itself. One key component of intelligence seems to be adaptiveness, the ability of the agent to change itself or its surroundings when it so desires. Mundane adaptiveness involves, for example, the agent updating its beliefs to reflect the effects of its actions, to accomodate new information, or to adopt a stand on some question. Similarly, the agent might adapt by changing other mental structures, such as its desires and intentions (resolving inconsistencies, adopting stands, etc.) or its skills. In all these examples, the agent is acting as the designer of its new state, as its own programmer.

As its own programmer, the agent needs some way of guiding its adaptations, some way of stating its intentions about the relations of parts of its mental structure at single instants and over time so that it can modify itself to satisfy these intentions. But this is just the problem described above facing any programmer, one addressed by functional specifications relating parts of mental states to others at the same time, to others at other times, or to the mental state as a whole. The agent modifies itself so that its new state still satisfies (is a model of) the set of self-specifications. But since the machine is doing this revision itself with only its current state to work with, the machine's interpretation of these self-specifications cannot refer to external affairs, but must be narrowly interpreted specifications

referring only to parts of the agent itself. In classical terminology, artificial intelligence programs are not merely rule-obeying machines, but are rule-following machines, interpreters of the rules, and as such must be able to grasp the import of their rules.

Examples of such self-specifications are abundant in artificial intelligence programs. Perhaps the clearest example is that of the recently developed reason maintenance systems, also called truth maintenance, belief revision, or dependency network systems.<sup>33</sup> In a machine based on reason maintenance techniques, the fundamental sort of self-specification states that if the current state of mind contains certain components and lacks certain other components, then the current state of mind should also contain some further component. These specifications are termed *reasons* or *justifications*. The mental components related by reasons may be mental attitudes like beliefs, desire, etc.; descriptions, procedures, etc.; or whatever the psychology of the agent employs as building blocks for mental states. The agent follows these self-specifications by deriving its current mental state from the current set of reasons, using groundedness principles to construct a set of mental components satisfying all the reasons (i.e. by falsifying an antecedent or by including the conclusion). The agent's thinking and acting may change the state of mind by adding or subtracting new reasons to the current set, or by switching to another model of the current set. In each of these cases, the machine revises its current mental state by using the reasons as guides to what mental components should be adopted or abandoned.

Reasons are not the only sort of narrowly interpreted self-specification employed in artificial intelligence. The least complex self-specifications are those simply declaring the existence of some component of mental structure. Most of the declarations of the so-called "knowledge representation languages" take this form, asserting the existence of a belief (a statement in a database), of a desire (a goal statement), or of an intention (an item on an agenda). More complex self-specifications relate two or more components of states of mind, including reference or coreference relationships (procedural attachments and equalities).<sup>34</sup> The structure sharing, inheritance, or "virtual copy" relationships so common in representational systems are simply self-specifications stating that one description should be considered as having certain components if other descriptions have those components as well.<sup>35</sup> Likewise, MINSKY'S "K-lines" can be viewed as self-specifications stating that the state of mind should contain one subset of (K-node) components if it also contains the corresponding enabling (K-node) component.<sup>36</sup> In fact, much of what goes by the name of "self-description" in artificial intelligence is not merely descriptive but instead normative, and so properly viewed as self-specification rather than self-description.

---

<sup>33</sup>See [STALLMAN AND SUSSMAN 1977], [DOYLE 1979], [LONDON 1978], and [MCALLESTER 1980].

<sup>34</sup>See, for example, [WEYHRAUCH 1980] and [SUSSMAN AND STEELE 1980].

<sup>35</sup>See, for example, [FAHLMAN 1979].

<sup>36</sup>[MINSKY 1980]

### III. The concept of mind

#### *Convergence of the theories*

§31. We can now not help but see a convergence of ideas between philosophy of mind and artificial intelligence. Philosophy began with the idea that psychological theories describe the mind in ecological terms, but abandoned that view in favor of the idea of narrow psychological theories since ecological facts offer little predictive power about behavior, even when they are accurate descriptions of mental states. On the other hand, the theory of machines arrived somewhat earlier at the notion of effectiveness, a stronger notion entailing narrowness. Descriptions of machines (either abstract or programmed) were early on effective, hence narrow, theories of machines states and behavior. But effectiveness was so strong as to obscure matters. Effectiveness of machine description implies rather direct physical realizability: this fact lies at the heart of the practice of programming. Because of their privileged status in these direct realizations, programs became identified with the machines realizing them. Considerable effort was required to regain the perspective of abstract, general, and multiple theoretical specifications of machines, as opposed to unique or privileged machine descriptions. In computer science, the recovered perspective proved crucial in ecological specifications of machines in their environment of application. In artificial intelligence, the recovered perspective proved crucial in the design of adaptive machines reconstructing themselves by means of narrow self-specifications. We attribute narrowness rather than effectiveness to the self-specifications of artificial intelligence because often the sorts of specifications were introduced as ideal (but narrow) specifications only approximated in their interpretation by an accompanying effective algorithm. While philosophy came to narrowness seeking predictive power, artificial intelligence came to narrowness seeking adaptive power.

#### *Possible minds defined*

§32. The two paths to narrowness of psychological theories suggest the importance of the idea for cognitive science. It might seem that the motivation for adaptiveness subsumes the motivation of predictive power, since the rationality of a particular adaptation involves the agent's expectations about the effects of the changes. But in fact the motivations are separate, since many of the adaptation applications leading to the idea of narrow self-specification need not be deliberate or considered adaptations, but may be automatic reorganizations choosing some possible adaptation without regard to comparative advantages or disadvantages. These converging motivations mean that while any psychological theory, narrow or not, may be of interest in ecological studies of minds in their environments, only narrowly realized psychological theories need matter to cognitive science. We draw on the apparent significance of this idea to turn it around and say that *all* narrowly realized psychological theories matter to cognitive science, that the set of such realizations forms the set of possible minds.

§33. Defining the set of possible minds as the set of narrowly realized theories has several advantages for cognitive science, advantages of neutrality on several important questions. This neutrality permits the use of these questions as dimensions for classification rather than as presuppositions of the science. Specifically, the definition is neutral on questions of psychological ontology, complexity, effectiveness, and determinateness.

The definition is neutral on psychological ontology because as long as the theory has some realization, it may posit any sorts of mental entities it wants. For example, one can express narrow stimulus-response psychologies in terms of relations between sensors and effectors; attitudinal psychologies in terms of relationships between beliefs and desires (or whatever attitudes one chooses); Freudian psychologies in terms of ego, superego, id, energy, and flows; and even the monolithic states of finite state machines. Thus the first dimension of classification of minds is by the sorts of mental components, by psychological ontology.

The definition is neutral on psychological complexity for reasons similar to its neutrality on psychological ontology. As long as the theory is narrowly realizable, it is a possible psychology no matter how trivial or complex it is. Possible minds may be as simple as a soda machine, or as complex as LEV TOLSTOY. Thus the second dimension of classification of minds is by their structural complexity, by the variety of ways their components may be combined to form mental states.

The definition's neutrality on effectiveness follows since effectiveness is a stronger notion than narrowness. The range of possible minds includes both mechanically realized psychologies as well as physically realized psychologies which might not admit mechanized realizations. The definition allows for universes in which the notions of narrowness and effectiveness coincide, and for universes like our own in which narrowness subsumes effectiveness. Moreover, it allows for universes in which effectiveness is even more restricted than in our own. This neutrality opens a whole range of new questions for physical theorists. One might conjecture from the amazing coincidence of effective calculability and recursive computability in our universe that other universes might offer analogous coincidences of effectiveness with the degrees of recursive unsolvability or with the degrees of computational complexity.<sup>37</sup> What form would the physics of such universes take? Need they have different laws, or might they differ from ours only in the values of the fundamental constants?

Finally, the definition of possible minds is neutral on the question of determinism, admitting both theories in which a mental state may have at most one possible successor state and theories in which a mental state might be followed by any of several others. Deterministic and non-deterministic Turing machines are examples of these. Of course, even non-deterministic psychologies can have deterministic realizations, and the existence of non-deterministic realizations is a question of physics.

§34. Defining possible minds to be narrowly realized theories has a seemingly unavoidable and perhaps unwelcome consequence in addition to the previously discussed advantages. The neutralities on psychological ontology, complexity, effectiveness, and determinateness add up to a neutrality on realizations. According to the definition, any satisfiable theory is a possible psychology, and by taking the entire domain of the model to be an agent, any realized theory is a possible mind. This makes the idea of narrowness of theoretical realizations analogous to the idea of closed systems in physics, and with similar importance, the laws of physics being taken to describe all and only closed systems. The generality of the definition means that any object and an accurate internal description of it constitute a possible mind. For example, supposing the known laws of physics correct, the universe is a possible mind, with the laws of physics as its psychology. Similarly, the U.S. economy and a correct economic theory constitute a possible mind, as does *Hamlet* together with a correct descriptive analysis, as do the natural numbers together with PEANO'S axioms (assuming their correctness). If the aim of cognitive science is to classify all possible minds, then it includes as subdisciplines not just psychology and artificial intelligence, but also physics, theology, model theory, sociology, etc., etc., etc. If cognitive science is so all embracing, what endeavor is *not* cognitive science?

My own feelings are mixed about this problem. On the one hand, I find the seeming intellectual imperialism of this view distasteful, but on the other hand I can offer some perspectives from which it

---

<sup>37</sup>[ROGERS 1967], [AHO, HOPCROFT, AND ULLMAN 1974], [GAREY AND JOHNSON 1979]

seems less so, even natural. First, the generality of the definition does no real harm, for the point of the science is to introduce distinctions, and the first action of anyone interested in thinking would be to introduce distinctions of domains so as to reinstate the disciplines in their traditional fields. Indeed, most disciplines have ways of viewing the whole world from their perspective, but modesty and common sense keep them from overstepping their most fruitful bounds. Second, the generality seems unavoidable if one wants the field to encompass both the trivial mind (e.g. the soda machine) and physically realized but non-effectively realizable minds. Bizarre and pointless minds need not be very interesting to anyone, but the scope of the science must include them precisely so that terms like "bizarre" and "pointless" can be substantively applied. Group theory has its "monster" groups: cognitive science needs its "psychopaths." Third, and finally, there is a hidden, underlying rightness to the definition. All possible minds must incorporate all possible mental structures, of that there can be no doubt. Yet what are the mathematical structures of physics, the logical theories of model theory, the programs of computer science, the constitutions of governments, and the articles of the faiths if not ideas in the minds of men?



## IV. Conclusion

§35. Cognitive psychology and artificial intelligence currently dominate cognitive science, and with them come the dominant methodologies of experimental data collection and speculative construction. If the arguments of this paper are understood, practitioners of these fields will at least have been introduced and perhaps converted to a methodological viewpoint that begins with the set of all possible minds and proceeds by formulating distinctions and classifying minds. The new viewpoint does not reject, but subsumes, the previous viewpoints since experiment and speculation remain useful in the problems of formulation and classification.

Even if one remains attached to one's familiar methodology and hesitant before the general conception of mind introduced above, I think important practical benefits follow from approaching the study of mind in the proposed way. I am not familiar enough with the literature of cognitive psychology to promise cognitive psychologists these benefits, but I am familiar with artificial intelligence, and feel the intelligibility of its literature might be increased by adoption of the proposed viewpoint: not just intelligibility to outsiders, but to insiders as well. One of the commonly acknowledged problems in artificial intelligence is the difficulty of telling what someone else has done. Almost every researcher has his own vocabulary and viewpoint, and while the gist of papers is usually intelligible, the details present more difficulties because one person's omitted explanatory trivialities may be someone else's stumbling blocks. While the classificatory viewpoint cannot by itself reconcile vocabularies or world-views, consciousness of the task of classification might strengthen incentives for scholarly analysis of comparable works. Moreover, the classificatory viewpoint also makes it possible to state and infer what results are scientific results and what results are engineering advances. The scientific aim of artificial intelligence is a question of formulation and existence: whether or not there are interesting psychologies realizable in Turing-equivalent machines. The answer to this question will likely involve an actual realization. The details of this realization are not scientifically interesting, but instead are matters of engineering. This is not meant to belittle engineering problems, merely to distinguish two endeavors. The situation is similar to that in mathematics, where in the context of a particular theory, an individual proof is without mathematical significance except for its conclusion. Of course in a wider context the proof will have importance as an example of a method for discovering analogous proofs. Likewise, the details of an intelligent machine's construction will be answers to important engineering questions, perhaps useful in constructing realizations of other psychologies, but scientifically insignificant as far the answered existential question is concerned. This distinction between scientific and engineering questions allows clearer explanations of the problems and results reported in artificial intelligence. Is the purpose of a paper to show how to realize feature  $f$  of the psychology  $\psi$  currently being worked toward (a scientific advance)? Or is it to show ways in which feature  $f$  of the machine  $M$  might be utilized independently of any particular psychology (an engineering advance)? Is the paper's purpose to modify the target psychology in light of changing estimates of feasibility of realizations? Or is it to modify the underlying machine so as to enhance feasibility of realization? The first two of these questions serve to mitigate the apparent scientific irrelevance of papers on programming techniques and programming systems by viewing them as engineering advances rather than psychological theories. Likewise, the second two questions serve to clarify the revolutionary claims often made in the literature. Papers changing the target psychology or underlying machine change particular scientific aims. Papers inventing a new programming technique do not, but advance the progress of the scientific or engineering investigation of a particular set of scientific aims. Artificial intelligence may well be yet a field most concerned with problems of formulation, but let us at least make clear what is being formulated and studied, namely pairs of psychologies and machines.

§36. The proposed perspective on the aims of artificial intelligence can benefit the clarity and ease of construction of programs as well as the clarity and construction of their exegeses. The advantages of

initial specification and subsequent implementation are well-known from programming methodology in computer science, even if intermediate specifications and verifications of implementations rarely matter in artificial intelligence's problems of formulation. But the proposed perspective can have a far more substantial impact on the construction of artificial agents than simply as an explanatory aid. Earlier, we introduced the idea that intelligent machines can interpret narrow self-specifications to guide their adaptations. If we take this suggestion seriously, then designing machines to continually construct and interpret sets of self-specifications becomes a powerful aid to programming. In this methodology, intelligent machines are always designed to construct and update their mental states by interpreting their self-specifications. Reason maintenance systems were mentioned earlier as an example of this idea, but so are the constraint-based programming systems currently enjoying interest. Constraints are nothing more than self-specifications, and in programming systems based on them one simply describes the psychology of the program as a set of narrowly interpreted logical statements, which the programming system examines to maintain the current mental state. Of course, in artificial intelligence there must be not just narrow, but effective and feasible ways of interpreting the specifications, so that one repeatedly reformulates the specifications of the psychology as one's store of feasible interpretation techniques changes. For example, the extant constraint-based programming systems fantastically restrict the complexity of psychological specifications because they incorporate means for interpreting only a few of the very simplest sorts of specifications. But these restrictions are not essential, and future systems might allow incremental addition of interpretation techniques for specifications of more complex logical form.

The proposed methodology can be expressed as the thesis of self-interpretation: *Thinking is a process of narrow self-specification and self-interpretation.*<sup>38</sup> That is, mental actions are all described and realized by means of self-interpretive psychological theories which augment themselves with new self-specifications, purge themselves of unwanted self-specifications, or reinterpret their existing self-specifications. In this way mental actions are understood in terms of the psychology itself (e.g. reasons, beliefs, desires, etc.), rather than in terms of the realization of the psychology (e.g. CONS's, RPLACA's, etc.).

---

<sup>38</sup>This thesis is a descendant of the comparable thesis about reasoning proposed in [DOYLE 1979].

## References

- Aho, A. V., J. E. Hopcroft, and J. D. Ullman, 1974. *The Design and Analysis of Computer Algorithms*, Reading: Addison-Wesley.
- Backus, J., 1978. Can programming be liberated from the von Neumann style? A functional style and its algebra of programs, *C. A. C. M.* **21**, 613-641.
- Barendregt, H. P., 1981. *The Lambda Calculus: its syntax and semantics*, Amsterdam: North-Holland.
- Boring, E. G., 1950. *A History of Experimental Psychology* (2nd ed.), New York: Appleton-Century-Crofts.
- Bretano, F., 1874. *Psychologie vom empirischen Standpunkte*, Leipzig.
- Chomsky, N. A., 1965. *Aspects of Syntax*, Cambridge: MIT Press.
- Descartes, R., 1637. Discourse on the method of rightly conducting the reason and seeking for truth in the sciences, in *The Philosophical Works of Descartes* (E. S. Haldane and G. R. T. Ross, trans.), V. I, Cambridge: Cambridge University Press, 1931.
- Dennett, D. C., 1978. Intentional systems, *Brainstorms*, Montgomery, Vermont: Bradford Books, 3-22.
- Doyle, J., 1979. A truth maintenance system, *Artificial Intelligence* **12**, 231-272.
- Elgot, C. C., and A. Robinson, 1964. Random access stored program machines, *J. A. C. M.* **11**, 365-399.
- Fahlman, S. E., 1979. *NETL: A system for representing and using real world knowledge*, Cambridge: MIT Press.
- Feigl, H., 1958. The 'mental' and the 'physical,' in *Minnesota Studies in the Philosophy of Science* (H. Feigl, M. Scriven, and H. Grover, eds.) V. II, Minneapolis: University of Minnesota Press, 370-497.
- Fodor, J. A., 1968. *Psychological Explanation: An introduction to the philosophy of psychology*, New York: Random House.

- Fodor, J. A., 1975. *The Language of Thought*, New York: Crowell.
- Fodor, J. A., 1981. Methodological solipsism considered as a research strategy in cognitive psychology, *Representations: Philosophical essays on the foundations of cognitive science*, Cambridge: MIT Press, 225-253.
- Frege, G., 1884. *Die Grundlagen der Arithmetik*, Breslau: Verlag von Wilhelm Koenner, translated by J. L. Austin as *The Foundations of Arithmetic* (2nd revised ed.), Evanston: Northwestern University Press, 1963.
- Frege, G., 1892. Über Sinn und Bedeutung, *Zeitschrift für Philosophie und philosophische Kritik* 100, 25-50. Also, On sense and reference (M. Black, trans.), in *Translations from the Philosophical Writings of Gottlob Frege* (P. Geach and M. Black, eds.), Oxford: Basil Blackwell, 1952.
- Freud, S., 1895. Project for a scientific psychology, *Standard Edition of the Complete Psychological Works of Sigmund Freud* (J. Strachey, tr. and ed.), V. I, London: Hogarth Press, 1966, 281-397.
- Gandy, R., 1980. Church's thesis and principles for mechanisms, in *The Kleene Symposium* (J. Barwise, H. J. Keisler, and K. Kunen, eds.), Amsterdam: North-Holland, 123-148.
- Garey, M. R., and D. S. Johnson, 1979. *Computers and Intractability: a guide to the theory of NP-completeness*, San Francisco: W. H. Freeman.
- Goldstein, H. H., 1972. *The Computer from Pascal to von Neumann*, Princeton: Princeton University Press.
- Hilbert, D., 1900. Mathematical Problems, lecture delivered before the International Congress of Mathematicians at Paris (M. W. Newson, trans.), *Bull. A. M. S.* 8 (1902), 437-479.
- Hintikka, J., 1971. Semantics for propositional attitudes, in *Reference and Modality* (L. Linsky, ed.), 145-167.
- Kaplan, D., forthcoming. *Demonstratives*.
- London, P. E., 1978. Dependency networks as a representation for modelling in general problem solvers, Department of Computer Science, University of Maryland, TR-698.
- McAllester, D. A., 1980. An outlook on truth maintenance, MIT Artificial Intelligence Laboratory, Memo 551.

- Minsky, M., 1967. *Computation: Finite and infinite machines*, Englewood Cliffs: Prentice-Hall.
- Minsky, M., 1980. K-lines: a theory of memory, *Cognitive Science* 4, 117-133.
- Peters, R. S., and C. A. Mace, 1967. Psychology, in *Encyclopedia of Philosophy* (P. Edwards, ed.), New York: Macmillan and Free Press, V. 7, 1-27.
- Place, U. T., 1956. Is consciousness a brain process?, *British J. of Psychology* 47, 44-50.
- Plotkin, G., 1972. A set-theoretical definition of application, School of Artificial Intelligence, Memo MIP-R-95, University of Edinburgh.
- Putnam, H., 1975a. The meaning of 'meaning,' *Mind, Language, and Reality*, Cambridge: Cambridge University Press, 215-271.
- Putnam, H., 1975b. Minds and machines, *Mind, Language, and Reality*, Cambridge: Cambridge University Press, 362-385.
- Putnam, H., 1975c. Philosophy and our mental life, *Mind, Language, and Reality*, Cambridge: Cambridge University Press, 291-303.
- Quine, W. V., 1970. *Philosophy of Logic*, Englewood Cliffs: Prentice-Hall.
- Randell, B., (ed.), 1975. *The Origins of Digital Computers: Selected papers*, Berlin: Springer-Verlag.
- Rogers, H. Jr., 1967. *Theory of Recursive Functions and Effective Computability*, New York: McGraw-Hill.
- Ryle, G., 1949. *The Concept of Mind*, London: Hutchinson.
- Scott, D., 1973. Models for various type-free calculi, *Logic, Methodology and Philosophy of Science IV* (P. Suppes, L. Henkin, A. Joja, Gr. C. Moisil, eds.), Amsterdam: North-Holland.
- Shepherdson, J. C., and H. E. Sturgis, 1963. Computability of recursive functions, *J. A. C. M.* 10, 217-255. . .
- Skinner, B. F., 1957. *Verbal Behavior*, New York: Appleton-Century-Crofts.

- Smart, J. C. C., 1959. Sensations and brain processes, *Philosophical Review* **68**, 141-156.
- Stallman, R. M., and Sussman, G. J., 1977. Forward reasoning and dependency-directed backtracking in a system for computer-aided circuit analysis, *Artificial Intelligence* **9**, 135-196.
- Sussman, G. J., and G. L. Steele Jr., 1980. CONSTRAINTS—A language for expressing almost-hierarchical descriptions, *Artificial Intelligence* **14**, 1-39.
- Thomason, R. H., 1979. Some limitations to the psychological orientation in semantic theory, mimeo, University of Pittsburgh.
- Turing, A. M., 1936. On computable numbers with an application to the entscheidungsproblem, *Proc. London Math. Soc. Ser. 2* **42**, 230-265.
- Watson, J. B., 1914. *Behavior: An introduction to comparative psychology*, New York: Holt.
- Weyhrauch, R. W., 1980. Prolegomena to a theory of mechanized formal reasoning, *Artificial Intelligence* **13**, 133-170.
- Whitehead, A. N., and B. Russell, 1910. *Principia Mathematica*, Cambridge: Cambridge University Press.
- W. Wundt, 1874. *Grundzüge der physiologischen Psychologie*, Leipzig: W. Englemann. Also *Principles of Physiological Psychology* (E. B. Titchener, trans.), New York: Macmillan, 1904.