
EMERGENCE OF REPRESENTATION IN AUTONOMOUS AGENTS

MARK H. BICKHARD

**Cognitive Science, Lehigh University,
Bethlehem, Pennsylvania, USA**

A problem of action selection emerges in complex—and even not so complex—interactive agents: what to do next? The problem of action selection occurs equally for natural and for artificial agents—for any embodied agent. The obvious solution to this problem constitutes a form of representation, interactive representation, that is arguably the fundamental form of representation. More carefully, interactive representation satisfies a criterion for representation that no other model of representation in the literature can satisfy or even attempts to address: the possibility of system-detectable representational error. It also resolves and avoids myriad other problematics of representation and integrates or opens the door to many additional mental processes and phenomena, such as motivation.

AUTONOMOUS AGENTS AND INTERACTION SELECTION

In sufficiently simple interactive systems, there may be only one thing to do, and it is simply done all of the time. Slightly more complex would be the case in which multiple actions or interactions are possible but only one of them is triggered in any possible combination of system state and input. With only minimal interactive complexity, however, there will be conditions in which multiple interactions are possible, and a simple

Address correspondence to Mark H. Bickhard, Cognitive Science, 17 Memorial Drive East, Lehigh University, Bethlehem, PA 18015, USA. E-mail mhb0@lehigh.edu; <http://www.lehigh.edu/~mhb0/mhb0.html>

Thanks are due to the Henry R. Luce Foundation for support during the preparation of this paper.

triggering of one of those possibilities will be insufficient. This will occur, for example, if the system-environment circumstances have not occurred before, or have occurred too few times to have selected for a simple triggering relation, or when the environment is insufficiently predictable for it to be determinate what the system outcomes of the interactions will be (Bickhard & Terveen, 1995). A given frog tongue flick may or may not catch a fly.

Action selection in such cases poses its own special problems. In particular, whatever presuppositions are involved in the general strategy of interaction selection may, in particular cases, turn out to be wrong. In general, frogs flick tongues in order to catch flies (or worms . . .), but such outcomes are not guaranteed. In some cases, the advisable strategy would be simply to try again when the interaction does not work, whereas in other cases the advisable strategy might be to switch to a different next interaction selection.

But implementing any such design constraints requires that the system be able to detect when the “desired” outcomes have occurred or have not occurred. That, in turn, requires that those desired outcomes be somehow functionally indicatable and detectable. In other words, the agent has to be able to differentiate between interaction outcome success and interaction outcome failure, and this requires that it be able to be functionally sensitive to information about interaction outcomes.

There are many traps in the way to an implementation of these notions. For example, if the desired outcomes are certain *environmental* conditions—perhaps a new position relative to some food source—then those environmental conditions will have to be *represented* by the system, and failures to achieve those environmental conditions will have to be detectable. But representation is what we are aiming to model in the first place, so this approach initiates a fatal circularity.

There are multiple logical and conceptual dangers involved in this modeling, some of ancient provenance, some much more recently discovered (Bickhard, 1993; Bickhard & Terveen, 1995). I will not review these here but will move more directly to what I advocate as the solution: The relevant interaction outcomes must be *internal* so that they do *not* have to be represented but can be directly functionally accessible to the system. Such internal interaction outcomes must be *anticipated* (Bickhard, 1993; Brooks, 1994) in such a way that success and failure of those anticipations can be detected.

There are many architectures within which the functions of anticipation and success/failure detection can be implemented (Bickhard & Terveen, 1995). For my current purposes, I will focus only on simple and familiar architectural principles of states and pointers. In particular, interaction outcomes will be system *states*, so anticipations of system states are what is required. Anticipations, in turn, require only *indications* of such system states associated with the relevant *interactions* that might or should yield those states. And, finally, such indications can be accomplished with simple *pointers*. Interaction possibilities can be indicated with pointers to the subsystems or system conditions that would engage in those interactions if initiated, and associated outcomes can be indicated with associated pointers to the relevant internal states. As long as the system can detect whether or not those indicated states have been achieved, all of the functional requirements for this form of interaction selection have been met.

Selecting interactions on the basis of the indicated outcomes of those interactions permits many powerful phenomena. It permits strategic selection of interactions *toward* particular outcomes—*goal* directedness. It permits *error-guided strategies* of interaction. It permits error guidance of representational *learning*, as long as the system is capable of learning at all. Any complex autonomous agent will require such forms of interaction selection. Note, further, that interaction selection is the fundamental problem of motivation (Bickhard, in press-b; Bickhard & Terveen, 1995; Cherian & Troxell, 1995a; Maes, 1990a, 1990b, 1991, 1992, 1994).

REPRESENTATION

Autonomous agents encounter problems of interaction selection. A natural solution to those problems involves internal indications of internal interaction outcomes. The central claim of this paper is that such indications constitute a form of *representation*, arguably the fundamental form.

I will support this claim in three ways. The first two supports address two properties of representations that this model accounts for: (1) content and (2) truth value of that content. The third points out that the interactive model provides an account for a fundamental criterion for original representation that is not and cannot be addressed by any

standard model in the contemporary literature: the possibility of *system-detectable* representational error.

Representational Content

The claim that indications of interaction outcomes can constitute representation poses a number of questions. A central question is: What constitutes representational content in this model? The answer has, in fact, already been adumbrated above: an indication of an interaction and its associated outcomes makes presuppositions about the environment. In particular, those presuppositions will hold in some circumstances and not hold in others, and, correspondingly, the indicated outcomes will occur in some—the presupposed—circumstances and not occur in others. Such indications involve implicit presuppositions of the “appropriate” properties of the environment for those indicated outcomes to occur.

These presuppositions provide a model for representational content. The presupposed interactive properties of the environment are what are represented of that environment. They are what are *predicated* of that environment: “This environment has the interactive properties appropriate to these indicated interaction outcomes.” The content in this model is *implicit* rather than explicit (Bickhard, 1993; Bickhard & Terveen, 1995; Brooks, 1991b), presupposed rather than portrayed. This is a major strength of the model in that it avoids many aporetic problems concerning representation (Bickhard, 1993; Bickhard & Terveen, 1995). It also poses its own problems—for example, how then do we account for clear cases of explicit content? How can representations of objects be modeled, not just representations of interactive properties? How about representations of abstractions, such as numbers? Such questions lead into further elaborations and explorations of the interactive model. I have addressed many of these elsewhere (Bickhard, 1980, 1993, in press-a, b; Bickhard & Richie, 1983; Bickhard & Terveen, 1995). At this point, I simply want to point out that the interactive model does provide an account of content. The epistemological world is constructed out of such content (Bickhard, 1993; Bickhard & Terveen, 1995; Hooker & Christensen, in preparation; Kuipers, 1988; Kuipers & Byun, 1991; Matarić, 1991; Nehmzow & Smithers, 1991, 1992).

Truth Value

Furthermore, those presuppositions—those implicit contents—can be *false*, in which case the outcomes will not occur. Or they can be true, in which case the outcomes will occur. That is, the actual environment is the truthmaker for the implicit content, and, consequently, the interactive model of representational content can account for truth value, one of the fundamental properties of representation.

SYSTEM-DETECTABLE REPRESENTATIONAL ERROR

Still further, interactive representational true value can be detected, at least fallibly, by the interactive system itself. The outcome indications can be falsified by and for the system, when it detects that the indicated outcome did not occur. That constitutes a falsification of the presupposed properties, the implicit content, about the environment in which the interaction occurred. It is a falsification that is accomplished by the system itself.

Note that, unless such falsifications can occur by and for the agent itself, *error guidance* is not possible, whether interactive error guidance or learning error guidance. Error guidance, in turn, is required in any complex interactive system because the certitude required for, say, simple triggering is not possible.

System-detectable error, however, is not possible for standard models of representation in the literature. Standard models have difficulty accounting for any kind of error and do not even address system-detectable error. Error detection for Fodor (1987, 1990a, 1990b), Dretske (1981, 1988), and Millikan (1984, 1993), for example, requires comparisons between what actually triggers an input state and what “ought” to trigger that state. What ought to have triggered such a state, in turn, is determined by considerations of evolutionary histories and counterfactual asymmetric dependencies between classes of possible such triggers. The necessary comparisons, then, are possible, if at all, only for an observer of the system, external to the system, that might be in a position to determine the relevant histories and dependencies. In particular, they cannot be determined by and for the system itself. The frog knows nothing about the evolutionary history of its fly detectors, nor about any dependencies, asymmetric or otherwise, between fly detections and BB detections. Standard models cannot account for system-

detectable representational error; standard models cannot account for error-guided interaction or learning; standard models cannot account for representation for the agent itself. The interactive model can. [Furthermore, interactive representational content can “ground” internal symbols as stand-ins for those contents (Bickhard & Terveen, 1995).]

Fodor, Dretske, and Millikan at least attempt to address the issue of representational error, even if not *system-detectable* representational error. The infamous symbol system hypothesis does not. Representation of the environment, according to the symbol system hypothesis, is constituted as structures in the system that are isomorphic with whatever is to be represented in the world (Newell, 1980; Newell & Simon, 1975, 1987; Vera & Simon, 1993, 1994). But isomorphism is ubiquitous throughout the universe (e.g., in lawful causal relationships, among others); representation is not (Smith, 1987). Isomorphism is a reflexive relation; representation is not. Isomorphism is a symmetric relation; representation is not. Isomorphism is a transitive relation; representation is not. Isomorphism either *exists*, in which case the purported representation is correct, or it does *not* exist, in which case the purported *representation* does not exist, and, therefore, *cannot be incorrect*; along with all of its other problems, isomorphism cannot account for representational error (Bickhard, 1993; Bickhard & Terveen, 1995). Isomorphism has insufficient resources to constitute representation: it either exists or not and is therefore incapable of differentiating the third representational possibility of existing but being in error (Millikan, 1984). It is little wonder that such approaches to representation have encountered a thicket of unresolved issues (Bickhard, 1993; Bickhard & Terveen, 1995; Smith, 1987) and a reaction against the very notion of representation (Beer, 1990; Brooks, 1991a; van Gelder, 1995).

Connectionist notions of representation fare no better (McClelland & Rumelhart 1986; Rumelhart & McClelland, 1986). A trained net establishes a correspondence between its output state vector and some class of inputs and thereby differentiates the current input into that class (Bickhard & Terveen, 1995). But such correspondence is simply a degenerate version of the symbol system hypothesis of “representation as isomorphism,” and it suffers from all of the same fatal defects. Correspondence is simply an unstructured, point-to-point, isomorphism. Neither the trained correspondences of connectionism nor the (causal) correspondences of transduction (Bickhard & Richie, 1983; Fodor, 1986, 1990a, 1990b, 1991; Fodor & Pylyshyn, 1981) can account for the possibility of representational error, and they do not even attempt

system-detectable representational error (Bickhard, 1993; Bickhard & Terveen, 1995).

Note that these standard models all look externally and backward in time, whereas the interactive model looks internally and forward in time (Bickhard & Terveen, 1995). It is such backward-looking “spectator” models that arise naturally in symbol manipulation, information theoretic, or connectionist models, but such an approach has never been able to solve the basic problems of representation. It has, rightly, yielded a reaction against the reality or usefulness of the very notion of representation (Beer, 1990, 1995; Brooks, 1991a, 1991c; Prem, 1995). Interactive representation constitutes an alternative.

REPRESENTATION AND DESIGN

Interactive representation already exists. It can be found in relatively simple organisms (not to mention complex organisms) and in some contemporary robots (Brooks, 1994; Cherian & Troxell, 1995a, 1995b—with a major caveat concerning the emergent constitution of *function* in robotic agents: Bickhard, 1993; Bickhard & Campbell, 1997; Millikan, 1984, 1993). It solves a basic design problem and provides guidance for more sophisticated versions of such problems and their solutions; it solves and dissolves many philosophical problems (Bickhard, 1993; Bickhard & Terveen, 1995). Conversely, the study of robots and autonomous embodied agents *inherently* encounters the fundamental problems of representation and epistemology, whereas passive system approaches, whether computational or connectionist, do not (Bickhard, 1982; Bickhard & Terveen, 1995). Interactive dynamic systems—embodied agents—are the correct locus for the study and design of mentality (Beer, 1995; Bickhard & Terveen, 1995; Hooker et al., 1992; Malcolm et al., 1989; Port & van Gelder, 1995; Steels, 1994).

It also, however, initiates a progression through many additional mental properties and thereby guides further understanding of those properties both in natural systems and in the design of artificial systems. These include, for example, motivation, perception, learning, emotions, consciousness and reflexive consciousness, language, and rationality (Bickhard, 1980, 1993, forthcoming-a, forthcoming-b; Bickhard & Richie, 1983; Bickhard & Terveen, 1995; Hooker, 1995; Kinsbourne, 1988). Interactive representation constitutes the interface of the emergence of mental properties in functional interactive system properties.

REFERENCES

- Beer, R. D. 1990. *Intelligence as adaptive behavior*. San Diego: Academic Press.
- Beer, R. D. 1995. Computational and dynamical languages for autonomous agents. In *Mind as motion: Dynamics, behavior, and cognition*, ed. R. Port and T. J. van Gelder, 121–147. Cambridge, MA: MIT Press.
- Bickhard, M. H. 1980. *Cognition, convection, and communication*. New York: Praeger.
- Bickhard, M. H. 1982. Automata theory, artificial intelligence, and genetic epistemology. *Rev. Int. Philos.* 36:549–566.
- Bickhard, M. H. 1993. Representational content in humans and machines. *J. Exp. Theor. Artif. Intell.* 5:285–333.
- Bickhard, M. H. (in press-a). Critical principles: On the negative side of rationality. In *Beyond ruling reason: Non-formal approaches to rationality*, ed. W. Herfel and C. A. Hooker.
- Bickhard, M. H. (in press-b). Is cognition an autonomous subsystem? In *Computation, cognition, and consciousness*, ed. S. O’Nuallain. Amsterdam: John Benjamins.
- Bickhard, M. H., and D. T. Campbell. 1997. Emergence (manuscript).
- Bickhard, M. H., and D. M. Richie. 1983. *On the nature of representation: A case study of James J. Gibson’s theory of perception*. New York: Praeger.
- Bickhard, M. H., and L. Terveen. 1995. *Foundational issues in artificial intelligence and cognitive science—Impasse and solution*. Amsterdam: Elsevier Scientific.
- Brooks, R. A. 1991a. Intelligence without representation. *Artif. Intell.* 47(1–3):139–159.
- Brooks, R. A. 1991b. Challenges for complete creature architectures. In *From animals to animats*, ed. J.-A. Meyer and S. W. Wilson, 434–443. Cambridge, MA: MIT Press.
- Brooks, R. A. 1991c. New approaches to robotics. *Science* 253:1227–1232.
- Brooks, R. A. 1994. Session on building cognition. Conference on the Role of Dynamics and Representation in Adaptive Behaviour and Cognition, University of the Basque Country, San Sebastian, Spain, December 9, 1994.
- Cherian, S., and W. O. Troxell. 1995a. Intelligent behavior in machines emerging from a collection of interactive control structures. *Comput. Intell.* 11:565–592.
- Cherian, S., and W. O. Troxell. 1995b. Interactivism: A functional model of representation for behavior-based systems. In *Advances in artificial life: Proceedings of the Third European Conference on Artificial Life*, Granada, Spain, ed. F. Morán, A. Moreno, J. J. Merelo, and P. Chacón, 691–703. Berlin: Springer.
- Dretske, F. I. 1981. *Knowledge and the flow of information*. Cambridge, MA: MIT Press.

- Dretske, F. I. 1988. *Explaining behavior*. Cambridge, MA: MIT Press.
- Fodor, J. A. 1986. Why paramecia don't have mental representations. In *Midwest Studies in Philosophy X: Studies in the philosophy of mind*, ed. P. A. French, T. E. Uehling, and H. K. Wettstein, 3–23. Minneapolis, MN: U. of Minnesota.
- Fodor, J. A. 1987. *Psychosemantics*. Cambridge, MA: MIT Press.
- Fodor, J. A. 1990a. *A theory of content*. Cambridge, MA: MIT Press.
- Fodor, J. A. 1990b. Information and representation. In *Information, language, and cognition*, ed. P. P. Hanson, 175–190. Vancouver: University of British Columbia Press.
- Fodor, J. A. 1991. Replies. In B. Loewer, G. Rey (Eds.) *Meaning in Mind: Fodor and his critics*. (255–319). Oxford: Blackwell.
- Fodor, J. A., and Z. Pylyshyn, 1981. How direct is visual perception?: Some reflections on Gibson's ecological approach. *Cognition*, 9, 139–196.
- Hooker, C. A. 1995. Reason, Regulation, and Realism: Towards a Regulatory Systems Theory of Reason and Evolutionary Epistemology. SUNY.
- Hooker, C. A., and W. Christensen (in preparation). Very Simple Minds.
- Hooker, C. A., H. B. Penfold, and R. J. Evans. 1992. Towards a theory of cognition under a new control paradigm. *Topoi* 11:71–88.
- Kinsbourne, M. 1988. Integrated Field Theory of Consciousness. In *Consciousness in contemporary science*, ed. A. J. Marcel and E. Bisiach, 239–256. New York: Oxford University Press.
- Kuipers, B. J. 1988. The TOUR model: A theoretical definition. From Kuipers, B. J., and T. Levitt, Navigation and mapping in large-scale space *AI Mag.* 9(2):25–43.
- Kuipers, B. J., and Y. Byun. 1991. A robot exploration and mapping strategy based on a semantic hierarchy of spatial representations. *Robotics and Autonomous Syst.* 9:47–63.
- Maes, P. 1990a. *Designing autonomous agents*. Cambridge, MA: MIT Press.
- Maes, P. 1990b. Situated agents can have goals. In *Designing autonomous agents*, ed. P. Maes, 49–70. Cambridge, MA: MIT Press.
- Maes, P. 1991. A bottom-up mechanism for behavior selection in an artificial creature. In *From animals to animats*, ed. J.-A. Meyer and S. W. Wilson, 238–246. Cambridge, MA: MIT Press.
- Maes, P. 1992. Learning behavior networks from experience. In *Toward a practice of autonomous systems*, ed. F. J. Varela and P. Bourguine, 48–57. Cambridge, MA: MIT Press.
- Maes, P. 1994. Modeling adaptive autonomous agents. *Artif. Life* 1:135–162.
- Malcolm, C. A., T. Smithers, and J. Hallam. 1989. An emerging paradigm in robot architecture. In *Proceedings of the Second Intelligent Autonomous Systems Conference*, ed. T. Kanade, F. C. A. Groen, and L. O. Hertzberger, 284–293. Amsterdam, December 11–14, 1989. Stichting International Congress of Intelligent Autonomous Systems.

- Matarić, M. J. 1991. Navigating with a rat brain: A neurobiologically-inspired model for robot spatial representation. In *From animals to animats*, ed. J.-A. Meyer and S. W. Wilson, 169–175. Cambridge, MA: MIT Press.
- McClelland, J. L., and D. E. Rumelhart. 1986. *Parallel distributed processing*. Vol. 2: *Psychological and biological models*. Cambridge, MA: MIT Press.
- Millikan, R. G. 1984. *Language, thought, and other biological categories*. Cambridge, MA: MIT Press.
- Millikan, R. G. 1993. *White Queen psychology and other essays for Alice*. Cambridge, MA: MIT Press.
- Nehmzow, U., and T. Smithers. 1991. Mapbuilding using self-organizing networks in “really useful robots.” In *From animals to animats*, ed. J.-A. Meyer and S. W. Wilson, 152–159. Cambridge, MA: MIT Press.
- Nehmzow, U., and T. Smithers. 1992. Using motor actions for location recognition. In *Toward a practice of autonomous systems*, ed. F. J. Varela and P. Bourguine, 96–104. Cambridge, MA: MIT Press.
- Newell, A. 1980. Physical symbol systems. *Cogn. Sci.* 4:135–183.
- Newell, A., and H. A. Simon. 1975. Computer science as empirical inquiry: Symbols and search. In (1987) *ACM Turing Award Lectures: The first twenty years, 1966–1985*, 287–313. New York: ACM Press; Reading, MA: Addison-Wesley.
- Newell, A., and H. A. Simon. 1987. Postscript: Reflections on the Tenth Turing Award Lecture: Computer Science as empirical inquiry—symbols and search. In *ACM Turing Award Lectures: The first twenty years, 1966–1985*, 314–317. New York: ACM Press; Reading, MA: Addison-Wesley.
- Port, R., and T. J. van Gelder. 1995. *Mind as motion: Dynamics, behavior, and cognition*. Cambridge, MA: MIT Press.
- Prem, E. 1995. Grounding and the entailment structure in robots and artificial life. In *Advances in artificial life: Proceedings of the Third European Conference on Artificial Life*, Granada, Spain, ed. F. Morán, A. Moreno, J. J. Merelo, and P. Chacón, 39–51. Berlin: Springer.
- Rumelhart, D. E., and J. L. McClelland. 1986. *Parallel distributed processing*, Vol. 1: *Foundations*. Cambridge, MA: MIT Press.
- Smith, B. C. 1987. The Correspondence Continuum. CSLI-87-71, Center for the Study of Language and Information, Stanford, CA.
- Steels, L. 1994. The artificial life roots of artificial intelligence. *Artif. Life* 1(1):75–110.
- van Gelder, T. J. 1995. What might cognition be, if not computation? *J. Philos.* 92:345–381.
- Vera, A. H., and H. A. Simon. 1993. Situated action: A symbolic interpretation. *Cogn. Sci.* 17:7–48.
- Vera, A. H., and H. A. Simon. 1994. Reply to Touretzky and Pomerleau: Reconstructing physical symbol systems. *Cogn. Sci.* 18:355–360.