

Shared Mental Models: Ideologies and Institutions

Arthur T. Denzau and Douglass C. North¹

Center for Politics and Economics

Claremont Graduate School

and

Center for the Study of Political Economy

Washington University (St. Louis)

The rational choice framework assumes that individuals know what is in their self interest and make choices accordingly. Do they? When they go to the supermarket (in a developed country with a market economy) arguably they do act accordingly. In such settings, the individual knows, almost certainly, whether the choice would be beneficial, ex post. Indeed financial markets in the developed market economies (usually) possess the essential characteristics consistent with substantive rationality. However, it is simply not possible to make sense out of the diverse performance of economies and polities both historically and contemporaneously if individuals really knew their self interest and acted accordingly. Instead people act in part upon the basis of myths, dogmas, ideologies and "half-baked" theories.

We argue here both that ideas matter, and that the way that ideas are communicated among people is crucial to building useful theories that will enable us to deal with strong uncertainty problems at the individual level.² For most of the interesting issues in political and economic markets uncertainty, not risk, characterizes choice-making. Under conditions of uncertainty, individuals' interpretation of their environment will reflect the learning that they have undergone. Individuals with common cultural backgrounds and experiences will share reasonably convergent mental models, ideologies and institutions and individuals with different learning experiences (both cultural and environmental) will have different theories (models, ideologies) to interpret that environment. Moreover the information feedback from their choices is not sufficient to lead to convergence of competing interpretations of reality. In consequence, as Frank Hahn has pointed out, "there is a continuum of theories that agents can hold and act on without ever encountering events which lead them to change their theories" (Hahn, 1987, p. 324). In such cases, multiple equilibria will result. It is the argument of this essay that in order to understand decision making under such conditions of uncertainty we must understand the relationship between the mental models that individuals construct to make sense out of the world around them, the ideologies that evolve from such constructions, and the institutions that develop in a society to order interpersonal relationships. Let us begin by defining each concept.

Following Holland et al. (1986, p. 12), we start with the presumption that "...cognitive systems construct models of the problem space that are then mentally "run" or manipulated to produce expectations about the environment." For our purposes in this paper, ideologies are the shared framework of mental models that groups of individuals possess that provide both an interpretation of the environment and a prescription as to how that environment should be structured. As developed in North (1990, p. 3), institutions are the rules of the game of a society and consist of formal and informal constraints

¹The authors would like to thank John Nachbar and Randy Calvert, and seminar participants at the Washington University Economic History Lunch, Claremont Graduate School and the Public Choice Society. We apologize in advance for the remaining errors.

²The literature on finite automata starting with Aumann (1981) has taken a more formalist path to explore the implications of specific notions of bounded rationality (although some term this irrationality). This literature is compatible in certain ways with our argument, and can be supplemented by the communication of mental models notions developed here. For a comprehensive survey of this literature, see Binmore (1987, 1988) or Marks (1992).

constructed to order interpersonal relationships. The mental models are the internal representations that individual cognitive systems create to interpret the environment and the institutions are the external (to the mind) mechanisms individuals create to structure and order the environment.

Some types of mental models are shared intersubjectively. Different individuals with similar models enables them to better communicate and share their learning. Ideologies and institutions can then be viewed as classes of shared mental models. Our analysis in this paper is aimed at describing the more general set of shared models. The large work on cognitive science, especially the recent explosion of work on connectionism, can be used to analyze the features and dynamics of mental models, and thus of ideologies and institutions as well. But the social aspects of these models are of crucial importance in human society, and these cultural links are only now being explored in this literature (Hutchins and Hazlehurst, 1992). These social features are modeled in this paper as necessitating communication that allows an individual's experiential learning to be based on a culturally provided set of categories and priors so that each person does not need to begin as a *tabula rosa*.

The mental models that the mind creates and the institutions that individuals create are both an essential part of the way human beings structure their environment in their interactions with it. An understanding of how they evolve and the relationship between them is the single most important step that research in the social sciences can make to replace the black box of the "rationality" assumption used in economics and rational choice models. We need to develop a framework that will enable us to understand and model the shared mental models that guide choices and shape the evolution of political-economic systems and societies. What follows is an outline of how to go about this task.

I. The Chooser Facing Uncertainty and the Conditions for Substantive Rationality

Neoclassical economics has evolved, especially since Marshall left the scene, into a series of applications of the constrained optimization model, under complete information. Von Neumann and Morgenstern, followed by Savage and others in the 1940s and early 1950s, extended the model to incomplete information, so that the chooser faces risk and chooses a lottery rather than a unique outcome. However, the overarching presumption is that the resulting choices always reflect substantive rationality. This approach has been under attack by a few economists, as well as other social scientists and philosophers, from many viewpoints for decades, but there has been a lack of a serious alternative that incorporates the successful applications of the substantive rationality optimization model while still dealing in some productive manner with the shortcomings.

Friedman (1953, pp. 19-23) provides one of the fundamental defenses for the substantive rationality, as well as laying out its basic features (p. 21): he considers "the economic hypothesis that under a wide range of circumstances individual firms behave as if they were seeking rationally to maximize their expected returns...and had full knowledge of the data needed to succeed in this attempt; as if, that is, they knew their relevant cost and demand functions." Friedman states that it is unnecessary for the substantive rationality model to be a descriptive model, with the detailed implication true at the individual level. Rather, the model is supposed only to be applied empirically at the aggregated, or market, level. Even if we accept this justification and the philosophic approach behind it, there is still a problem for the substantive rationality paradigm: there are situations of societal decisions, or resource allocation, that substantive rationality models poorly predicts, even at the market level. Examples of this sort are provided in Section II. We believe that this is the situation now facing economics (and politics) in major areas of decisions, and that we must seriously consider the development of alternatives to applying substantive rationality to situations where it performs poorly.

The Conditions for Substantive Rationality

Before going further with this long march away from the neoclassical economist's behavioral assumption, we should further justify the need to take this divergent path. To

do this, we first consider a question little asked in the economic literature: What characterizes the domain of application of the substantive rationality paradigm?

One way to approach this problem is to consider a simple situation of choice in which substantive rationality models work well. Consider choice in competitive posted-price markets. In such a situation, the chooser need only choose the quantity to buy or sell, as the competitive environment makes the agent's situation relatively simple - the price can effectively be viewed as a parameter, and only the quantity need be chosen at this parametric price. The experimental literature, casual empiricism and much empirical literature (at least on the demand side) have shown this to be a good predictive model. Both proponents and critics typically acknowledge the power of the competitive behavior version of the substantive rationality paradigm in the appropriate domain of application. Even more widely studied, and found very successful has been the experimental study of Double Auctions (DAs).

However, recent work by Gode and Sunder (1992a, 1992b, 1992c) raises important questions about why the substantial rationality approach is so successful in the DA setting. They measure efficiency success by the percentage of sum of potential buyer and seller rents (also termed consumer surplus and profits) that are realized. The first and second papers calculate the expected efficiency for several different types of exchange institutions, using traders they term Zero Intelligence (ZI). These traders "lack power to observe, remember, search, maximize, or seek profits." They summarize each case by the minimum of the expected efficiency.

In a sealed-bid auction, they find that the minimum efficiency for unconstrained ZI traders (ZI-U) is 0. If the bidders are constrained to only make bids that do not yield them losses on the proposed trade, these constrained (ZI-C) traders generate a minimum efficiency of 75%. Gode and Sunder (1992a cited in 1992b, p.2) found that employing profit-motivated human traders instead of the ZI-C traders improved efficiency by 1%. They then compared their ZI-C traders when placed in two different versions of a standard experimental double auction. One auction allowed the bidders to accept bids and make contracts continuously, while the other mechanism first waited for all bidders to submit a bid before trying to clear any contracts. This difference in institution alone, with traders who do not respond at all to the differing strategic opportunities available, raises the minimum expected efficiency from 75% to 81%. The 6% improvement due to an institutional change, compared to the 1% improvement using human subjects, suggests that institutional features by themselves can be as important as rationality in generating efficient economic performance. The 81% minimum for one of the auction institutions even suggests that most efficiency gains in some resource allocation situations may be attributable to institutional details, independent of their effects on rational traders.³

Individual Chooser Attributes

What are the features of the choice environment that make the substantive rationality model work so well in the posted-price case? We believe that the following are the most important, but further study of the experimental literature may require updating:

Complexity - How complex are the mental models required in order to make sensible choices given one's preferences and resources? This can best be judged by the similarity of the most appropriate models (so far as we now know) to what the chooser already knows. Analysis does not reduce the complex to the simple. Rather, it's a process by which we substitute a familiar complexity for one that we have found novel. The

³This may be the explanation for the results first noticed by Ledyard over a decade ago, and presented in seminars at that time, and which have resulted in Easley and Ledyard (1992). They study the individual decisions by participants in experimental double oral auctions, as Ledyard had found that they were not behaving as predicted by economic theory; i.e., they had apparently not learned the "right" model, even though the allocation results were very efficient (over 95% of the consumer and producer rents captured by the parties).

invisible hand result is now obvious and intuitive not because it is simple, but because we are trained to see it when it may be present or useful. Thus, complexity and the frequency of similar choices may be related.

Motivation - This has two aspects. First, how important is the choice and the models underlying it to the individual? If the choice involves issues that are central to how the individual assesses themselves and the world, and is paying substantial attention to the situation, then learning may be much more rapid. Larger cognitive resources, both in terms of time and attention, are likely to be allocated to evaluating the choice and its effects. Second, to what extent does the individual believe that her own choice can affect the real outcomes? For example, the choice may be one that is actually made by some collective body such as a committee. In such a case, the individual may realize that she cannot control the decision and might devote less effort. On the other hand, the individual may perceive the situation as one that has carryover benefits to other such situations, and treats the learning as a capital investment with payoffs beyond the specific situation in which it is presented.

There is already experimental evidence of the existence of significant cognitive costs in learning to make decisions, and of its capital investment nature. Jamal and Sunder (1988) have found when untrained subjects are involved in an auction experiment without being paid dollars for their accrued buyer and seller rents, the convergence to a competitive equilibrium price is slow and there is substantial variance even at the end of an auction period. Note that the efficiency achieved does not differ much whether the subjects are already trained, or not, but the Gode and Sunder work apparently provides an explanation for the high efficiency in the untrained auction experiments. When the same class of subjects are placed in the same institution, but know they are going to be paid, the convergence is much faster and variance much lower. In fact, once these subjects have gone through a paid session, and then are in an unpaid auction experiment with the same rules, the results are the same as if they were being paid. This strongly supports the idea that learning is an irreversible investment.

Information (quality and frequency) - How good is the information provided that would allow one to correct bad models? That feedback is essential to learning is suggested by the Jamal and Sunder experiments. The feedback needs to be in a form that makes its relevance to the mental models transparent, or complexity is increased further.

How often does the choice occur, or similar choices, in a situation in which feedback is provided? This problem is crucial in complex models. The basic problem is that the mappings we are trying to learn are usually multidimensional, possibly involving several dimensions in a complex, nonlinear relation. We only have a finite, often very small, data sample of real experiences from which to learn this mapping. This is not a simple statistical problem, especially when we start out not certain as to the relevant arguments involved in the mapping.

Easy Choices - Competitive Markets

Complexity - Minimal modeling is required for less complex situations. The institutions themselves may help reduce the complexity of the mental models that one must attempt to create and learn the parameters of the models. For example, behaving parametrically with respect to the market does not require building models of other agents. One need only decide how to maximize one's own utility. Coursey and Mason (1987) found experimentally that people can maximize unknown functions in a few (5 to 10) choices if they are told the value of the function after each choice proposed.

Information - Having good information is essential to improving the mental models. Most of the information needed to make decisions is readily available in posted price settings, and, due to the frequency of choices, there is a large enough sample to improve estimation of the necessary empirical relations. The consumer in a competitive market makes similar choices continually, and may be buying the same items continually. Even if the product is not being bought frequently, the choice is like many others, and protocols for dealing with such occasional choices have been developed and refined.

Motivation - Feedback is direct in some private goods markets; for search goods, the feedback is immediate. Even for experience goods, the feedback is a bit less immediate, but still reasonably quick. This holds nicely for private goods, but becomes more problematic for non-private goods, as suggested by Down's idea of rational ignorance regarding publicly provided bundles of goods.

II. Strong Uncertainty, and More Complex Problems

Strong, or Knightian, uncertainty would occur when a chooser cannot be viewed as capable of having even subjective probability distribution functions defined over a set of possible outcomes. Likely cases occur when the chooser cannot even state a list of outcomes that the chooser would distinguish in terms of their values. Without such a list, one cannot act as though the situation is one of Knightian risk or of Savage subjective probabilities. We believe that all people start out life in such a situation of strong uncertainty. Holland et al. argue that one needs to organize one's observations and learning into some sort of structure; one not already programmed at birth, and we discuss some of the implications of such an approach to knowledge representation in sections III, IV and V.

If all choices were simple, made frequently with substantial and rapid feedback, and involved substantial motivation, then substantive rationality would suffice for all purposes. It would be both a predictive and descriptive model of equilibrium settings, and learning models based upon it could be used to describe the dynamics out of equilibrium.

But not all choices have all these characteristics. One problem is that one is not even certain whether a particular choice will improve one's circumstances or not. The choice may be made infrequently, sometimes only once in a lifetime. Without direct experience, information about potential outcomes may not be known or easily acquired. In these circumstances, substantive rationality may not be a good descriptive model. In some of these cases, however, the Gode and Sunder results suggest that the substantive rationality models can still be predictive even if rationality is actually irrelevant to the human behaviors involved.

But there are hard choices, made in institutional settings that are not as conducive to efficiency as the double auction. It is these problems that are now coming to the forefront in the social sciences. We have already developed an adequate framework for the easier problems in which the substantive rationality gives good results. But we have (for the most part implicitly) sometimes made the erroneous assumption that we can extend without explicit consideration the scope of the substantive rationality assumption to deal with the problems of ambiguity and uncertainty that characterize most of the interesting issues in our research agenda and in public policy. Problems in political economy, economic development, economic history for example, all require an understanding of the mental models and ideologies that have guided choices. It is now time to refocus on the wide range of problems that we have so far ignored that involve strong uncertainty.

Let's consider a likely candidate for being a hard choice. Suppose you are making a choice as to accepting or rejecting a take-it-or-leave-it offer of 10 apples for \$3 in a situation isolated from other potential apple sellers. Suppose further that you believe that the apples are of such a quality that you value them at more than \$3. In order to increase your utility, however, you would like to acquire them at a lower price. You may wish to assess whether the seller is really willing to walk away if you reject the first offer, or would begin to bargain. You need to begin building a mental model of the seller, based on whatever information is available. Your past interactions, both with the seller and with other such vendors in similar situations, provide information for this construction. If you are purchasing 1,000 such bags, the motivation involved becomes more salient.

Further, the situation may, or may not, be one of potential continual dealing. You must be able to assess this probability in order to best evaluate the risk that the apples are not what you expect them to be. All of these factors require, as Arthur (1992) suggests, the building of internal representations of the agents with whom one interacts.

Even harder would be a game situation involving a multiplicity of other agents. In such a situation, the likelihood that substantive rationality holds begins to dwindle more rapidly. The complexity of the situation dramatically increases, as is pointed toward by Kreps, Milgrom, Roberts and Wilson (1982). Arthur (1992, p. 4) argues that the obvious flaws that would exist in one's mental models of other agents would make such decision situations not well-defined, and thus the game becomes a situation in which the standard rational choice framework has no application and no results.

Another dimension of hard choices is involved in collective choices. Downs (1957, ch. 14) argued that the reduced ability of an individual to determine the decision also reduces the incentives to become informed about it. We argue that this also reduces the motivation to allocate cognitive resources to the building and improving of mental models.

In order to deal with the variation in the complexity of decision problems, Arthur (1992, p. 5) has introduced the idea of a problem complexity boundary. In dealing with problems less complex (Lindgren, 1992, discusses the problem of a metric for this complexity) than this boundary, the substantive rationality approach is often a successful modeling approach, even if not all individuals would perform the problem analysis perfectly. The standard economic approach serves well with these problems. But with problems beyond the complexity boundary, Arthur argues effectively that the deductive rational procedure cannot be relied on. He claims that "the level at which they can do this [use the deductive rational procedure] reliably and accurately, I believe, is surprisingly modest." In spite of this, we are able to make decisions even in situations that are not well-defined, and thus in which the rational procedure provides no clue as to how to proceed.⁴

In these situations, we must be using some procedure that differs fundamentally from the deductive rational procedure. But what is that procedure? i.e., how can people make choices when faced with complex problems in a situation of strong uncertainty? Holland et al., and Arthur (1992) argue that we must be employing some form of induction, enabling us to learn from the outcomes of our previous choices. To usefully learn by induction, an individual needs some sort of mental model with which to understand the implications of a chosen action, as well as needing some way to identify potentially useful actions and the possible outcomes of those actions. The very spaces for actions, outcomes and reasonable strategies, as well as the mappings between them may be objects of ignorance on the part of the individual.

If problem complexity is too great, possibly caused by unreliable information as to the state of the world, then the substantive rationality results do not hold. Modeling such situations require one to model the decisionmaker as building internal mental models to represent the world and to learn from that world in order to improve the resulting choices.

III. Choice and Strong Uncertainty

Heiner (1983) presents a complementary argument in situations of uncertain choice. He argues that when there is a gap between an agent's competence and the difficulty of the decision problem to be solved (a C-D gap), the human agent constructs rules to restrict the flexibility of her own choices in such situations - i.e., institutions. This result can be derived using expected utility analysis once one incorporates a lack of reliability in interpreting environmental signals. By channelling choices into a smaller set of actions, an institution improves the ability to perceive the environment and to communicate. These benefits can then improve the ability of those involved in the institution to extract the potential gains from exchange or cooperation in production.

But that is not all that the agent does. Humans also construct explanations in the face of ambiguity and uncertainty and act upon them. In primitive societies we describe such explanations as myths, dogmas, taboos but in our own society we have religions,

⁴The finite automata approach is one way to deal with ill-defined problems, but this literature does not necessarily result in substantive rationality results.

superstitions, and other belief structures to account for many aspects of the environment for which we do not possess the information or acquire the feedback from choices made using that belief structure to arrive at something like a scientific consensus. How do we account for belief in such ideologies and act upon them when they entail faith?

A partial answer may be derived from an experiment by a psychologist, Julian Feldman (1959), in which subjects were shown sequences of 1's and 0's and were asked to predict which number would appear next. The subjects were quick to spot patterns in the sequence and to form hypotheses on the process generating the sequence. Using their models, they made predictions about which number would appear next when in fact the generation had been purely random.⁵

It may be an evolutionarily superior survival trait to have explanations for inexplicable phenomena, or this effect may just be a by-product of the curiosity which helps make humans model builders. But whatever is the explanation, mental models and ideologies play a crucial role in choice making.

IV. Learning and Shared Mental Models

In order to deal with the issue of how the mind copes with complexity we need first to step back and explore how learning occurs (Holland et al, 1986; Churchland, 1989; and Clark, 1989). There are two conceptually distinct levels at which learning can occur, with important implications for the effects of the learning. First, learning entails developing a structure by which to make sense out of the varied signals received by the senses. The initial architecture of the structure is genetic but its subsequent development is a result of the experiences of the individual. This architecture can be thought of as generating an event space which gets used to interpret the data provided by the world. The experiences can be classified into two kinds - from the physical environment and from the socio-cultural linguistic environment (Hutchins and Hazlehurst, 1992). The event space structure consists of categories - classifications that gradually evolve from earliest childhood on in order to organize our perceptions, and keep track of our memory of analytic results and experiences. Building on these categories we form mental models to explain and interpret the environment, typically in ways relevant to some goal(s) (Holland et al., p. 22). Both the categories and mental models will evolve to reflect the feedback derived from new experiences - feedback that may strengthen and confirm our initial categories and models or that may lead to modifications - in short, learning. Thus, the event space may be continually redefined with experience, including contact with other's ideas. Learning which preserves the categories and concepts intact, but which provides changed ideas about details and the applicability of the existing knowledge is the second level of learning. Together, learning within a given set of concepts and learning which changes the structure of concepts and mental models suggest a widely known approach to the dynamics of learning, which is further investigated in sections V and VI.⁶

⁵The Feldman experiment is discussed in Arthur (1992, pp. 12-3), whose intention was to show how human decision makers discern patterns in the context of complicated and ill-defined problems. In fact, what Feldman is showing is that individuals see patterns where they don't exist. In the Feldman experiment, as in life and science more generally, the models are underdetermined by the data. In other words, many models fit any finite data sequence, and data alone cannot judge between this multiplicity of "generalizations." Instead, one needs theory to generate hypotheses that can be tested, and impose constraints across sets of hypotheses involving different data in order to usefully perform inductions.

⁶Randy Calvert has suggested a means of formalizing these ideas. The action-outcome mappings can be defined as $o = g(a \mid s)$, where a is an action, s an equivalence class of situations representing the state of the world as viewed by the agent. The function maps into the outcome, o , or a pdf over outcomes. The agent also values this outcome with a utility function, $u(o \mid s)$. In order to avoid problems with the notation as the event space changes, we collapse this into a mapping from actions into utility, $f(u \mid a, s)$. The

It is at this juncture that the learning of humans will diverge from that of other animals (such as the sea slug which appears to be a favorite research subject of cognitive scientists) and certainly diverges from the computer analogy that dominated so much of early studies in artificial intelligence. The mind appears to order and reorder the mental models in successively more abstract form so that they become available to process information outside its special purpose origins. The term used by Clark and Karmiloff-Smith (forthcoming) to describe this process is representational redescription. The capacity to generalize, to reason from the particular to the general and to use analogy are all a part of this redescription process.

At the individual level, the representational redescription is a reorganization of the categories and concepts that is a form of learning distinct from the parameter updating that is occurring in the "normal learning" phase. Once a useful set of categories and concepts have been initially acquired, the normal learning period is long relative to the often sudden shifts in viewpoint that accompany representational redescriptions. The resulting dynamics are that of a punctuated equilibrium, first used by Eldredge and Gould (1972) to describe their new theory of biological speciation.⁷ Our application of the punctuated equilibrium idea is further delineated in section VI.A.

The world is too complex for a single individual to learn directly how it all works. The entire structure of the mental models is derived from the experiences of each individual - experiences that are specific to the local physical environment and the socio-cultural linguistic environment. "It follows that if two people have been exposed to different experiences in the past, with resulting differences in the stock of conceptual representations they have formed, they may act on the same data differently" (Arthur, 1992, p. 8). In fact, no two individuals have exactly the same experiences and accordingly each individual has to some degree unique perceptions of the world. Their mental models would tend to diverge for this reason if there were not ongoing communication with other individuals with a similar cultural background.

One of the crucial tasks of human development is to replace the nearly tabula rasa situation at birth with one informed extensively by various forms of indirect learning. The vast diversity of human culture that anthropologists have discovered suggests the relevance of this claim. In such a situation, learning other than direct must be providing the degree of similarity that one finds within each human society. The cultural heritage provides a means of reducing the divergence in the mental models that people in a society have and also constitutes a means for the intergenerational transfer of unifying perceptions. We may think of culture as encapsulating the experiences of past generations of any particular cultural group and, with the diversity of human experiences in different environments, there exists a wide variety of patterns of behavior and thought.

This learning can be called cultural learning, and what it provides in a pre-modern society is exactly the categories and concepts which enable a member of that society to organize their experiences and be able to communicate with others about them. Cultural learning in pre-modern societies not only provided a means of internal communication but also provided shared explanations for phenomena outside of the immediate experience of the members of the society in the form of religions, myths and dogmas. As noted earlier,

learning process is represented by the evolution of the equivalence classes over the situation space, as well as a Bayesian learning involving a fixed situation space, in which only the mapping between actions and utility change.

⁷If the event space does not undergo a representational redescription, then the dynamics are continuous. However, as footnote 4 notes, we are extremely unlikely to ever learn the true model, in the sense of assigning positive probability mass to it. It seems likely that there is always potential for learning more about the best event space in which to represent the world, and thus eventually the possibility for learning to result in a punctuation.

such belief structures are not confined to primitive societies but are an essential part of the belief structure of modern societies.

The rapid changes in lifestyles and technology of the past centuries has led to a proliferation and elaboration of ideologies. Each attempts to provide positive mental models that tend to focus on the actions and valued outcomes defined as crucial to hindering or fostering the vision embodied in the ideology (Downs, 1957; Higgs, 1987; Munger and Hinich, 1992). The positive mental models in an ideology comprise action-outcome mappings which relate the utility-relevant outcomes to the possible actions that the individual could choose among. For example, consider Milton Friedman's discussion about the social responsibility of business. He argues that the best way to be socially responsible, which we assume here to be an argument in the chooser's utility, is to maximize profits. Sowell (1980) further develops this argument, showing the crucial problems of information in attempting to deal with the effects of one's actions on unknown others.

Given the action-outcome mappings of an ideology, the normative or vision parts of an ideology identify the aspects of reality that are crucial to achieving one's goals. A Marxist would see the employment relation as an exploitive one: all profits produced in the capitalist production process results from the extraction of "surplus value" from the workers by the capitalist employing them, as the workers are induced to work for lower wages than the value of their labor. In attempting to examine the extraction of any excess value, a Marxist economist would attempt to measure the surplus value seized by the capitalist employer. A study of the strategies used to increase the surplus taken might then go on in order to determine what the workers' movement should spend its energies fighting, and what they should ignore. Using only this view of the world, one is likely to ignore many important changes that might make 1993 different from 1848, when Marx published the Communist Manifesto. We can attempt to see how the use of a shared mental model affects direct experiential learning.

IV.A. Learning in the Face of Strong Uncertainty with A Shared Mental Model

Let's see how to build a model of a chooser facing strong uncertainty, a chooser who learns both directly from the world and from a shared mental model (SMM). Figure 1 shows the basic framework of the uncertain chooser, learning directly from the external environment to improve his mental models. This process is slow, and can be made more rapid by having some indirect learning in the form of artifactual models already created by others, termed shared mental models.

This learner with shared mental models is shown in Figure 2. Here, the SMMs are related to the idea of Bayesian priors in a Bayesian learning model. But the Bayesian approach implicitly assumes that the dimensions of the internal mental models used to represent the external world are correct, in some sense.⁸ The connectionist approach and the classification models used by Holland et al. instead assume that the fundamental issue is to determine the relevant dimensions of reality for one's decision or learning purposes. For the learner, these dimensions are identified in large part by the existing shared mental

⁸This is one rationale for the result of Kalai and Lehrer (1990). They argue that unless the true model of the world is already given atomic mass (strictly positive probability assigned) in the support set of the learner's prior distribution, it is impossible for the learner to ever have the true model in the support set of any posterior distribution. If the attribute space in which the distributions are defined by a learner's mental models cannot be mapped in a straightforward way into a space in which the true model can be "naturally" located, it would be impossible to learn the true model using Bayesian methods.

models. A set of prior beliefs about action-outcome mappings is being learned as part of the shared mental model, whether traditional culture or ideology. ⁹

IV.B Mental Models in a Simple Model of Communication

The SMM has another important effect. It provides those who share the SMM, at least in the sense that they have an intellectual understanding of it, with a set of concepts and language which makes communication easier. Better communication links would lead to the evolution of linked individuals' mental models converging rather than diverging as they continue to learn directly from the world.

Figure 3 shows the idea of communication suggested by the Churchland view of knowledge representation. Agent L (for Local) has made a decision inside her mind, and wishes to explain the basis for the decision to her supervisor, agent C (for Center). The patterns in L's mind must first be encoded in a language, such as English. This encoding would be perfect if there were a known set of dimensions in which to measure the factors that caused L to make the choice she did, and if she could state her measurements of each of these dimension. This would constitute sufficient statistics for the decision, and communicating this data would be a perfect substitute for the neural patterns in L's mind.

But the problem is that we almost never know what the factors that result in a decision that we have made. Much of our understanding in a choice situation can be tacit knowledge, as Michael Polanyi discusses. We perceive things which we are not even consciously aware of, and this data can affect a decision. Attempts to determine the factors and their weights can be made, but the basic problem is that we are always uncertain as to the dimensions of the knowledge space that must be measured. As a result, the encoding is almost certainly to be imperfect, and not all the information used by L to make the decision can be placed in the communication channel.

The communication channel itself may be noisy and imperfect, and this problem has been studied extensively. This problem is a purely technical one, and is not the cause of the problems on which we wish to focus here. Instead, the decoding process at the listener, C, causes the next important communication problem. The listener must transform the message in the communication channel into changes in the neural patterns in his mind. The decoding is affected by the pre-existing patterns already in the listener's

9A second tie between our approach here and a Bayesian learning model needs to be laid out. The idea of a representational redescription is just like a surprise to a Bayesian learner: both seem to be impossible to generate from a Bayesian model. We believe, however, that this is a mistaken interpretation of the Bayesian model and not demanded by the model itself. Many of the punctuations in the learning of an individual result from the failures of a mental model to predict in situations when the individual is highly motivated, i.e., when the issue is one very important to the learner. The failure of the mental model to predict in such an important situation makes the person wish to avoid the negative reinforcement (opportunity cost that has been realized) in the future. Such motivation to learn causes the learner to mull over the problem to find its cause, and this process of reconsideration of the mental models can be viewed as giving substantial weight to the new data (the model's failure in a salient situation).

In a connectionist model, such repeated retraining on the new observation(s) can cause substantial changes in the connection weights, and thus the implicit concepts and relations embodied in the model. While this process has not been explicitly simulated in the experiments performed by cognitive scientists, this is because their mental models have not suggested the relevance of this idea. In a Bayesian model, the idea can be interpreted as the addition of new data with substantial (non-atomic) weight attached to it. Such a new mass of data different from the priors would cause a discrete jump in the posterior distribution from that prior. Both these types of discrete changes can be interpreted as the counterpart of representational redescriptions. These changes would also generate the type of punctuated equilibrium dynamics considered in section VI.A of this paper.

mind. The reception of a message and its interpretation by the listener are strongly influenced by the categories and beliefs that the listener already has about the world.

To the extent that the speaker and a listener have common features in their mental models for the concepts identified in the SMM, they are more likely to be able to encode and decode their internal ideas into a shared language, and more likely be able to effectively communicate using single terms to stand for substantial pieces of implicit analysis embodied in the SMM. To use the example of Marxian ideas again, consider the terms, "worker" and "capitalist." To use these words in the standard classical Marxian manner is to implicitly bring in considerable pieces of analysis of the exploitation of workers in a capitalist system. The terms also may carry affective denotations, so that the listener is expected to favor the worker and disfavor the capitalist. The world seen through this set of concepts can be a world quite different from that taught in a neoclassical price theory text.

By having a SMM available, the concepts embodied in the structure of mental models that several people have can be made more similar. As noted, the words used to convey the mental model ideas are used repeatedly as the espousers discuss their ideas among themselves, either orally or in written form. This should make the mental models of the two people more similar, and should enable their learning from some data observed by one of them relatively similar, compared to a random pair of individuals.

The sharing of mental models is enabled by communication, and communication allows the creation of ideologies and institutions in a co-evolutionary process. The creation of ideologies and institutions are important for economic performance, as there exist gains from trade and production that require coordination. As various authors have written, a market economy is based on the existence of a set of shared values such that trust can exist. The morality of a businessperson is a crucial intangible asset of a market economy, and its nonexistence substantially raises transaction costs. La Croix (1989) develops a model in which this intangible asset becomes a group-specific asset for a homogeneous middleman group (such as Jewish, Indian or Chinese traders in a society in which they are a minority). A small group that maintains itself differentiated from the rest of society can enjoy much lower transaction costs than would be true between two randomly chosen members of the society, and enable more transactions than would occur otherwise.

V. Mental Models and Institutions

How do mental models, institutions and ideologies interact to shape choices and the outcomes that determine political and economic performance? Mental models, institutions and ideologies are all a part of the process by which human beings interpret and order the environment. Mental models are, to some degree, unique to each individual. Ideologies and institutions are created and provide more closely shared perceptions and ordering of the environment. The connection between mental models and both ideologies and institutions crucially depends on the product and process of representational redescription. Both are, at this stage in cognitive science, quite imperfectly understood. The process involves an understanding of exactly how the progression in human cognition occurs. We believe that the punctuated equilibrium concept can be used to formalize this type of dynamics, as is discussed in section VI.A.¹⁰

The product has been more extensively analyzed than the process since there is a substantial psychology literature detailing experiments in learning. Here the agent's rate of learning varies with the difficulty of discerning expected pay-offs and "human learning can lock in to an inferior choice , and that this is prone to happen where pay-offs to choices are closely clustered, random and difficult to discriminate among" (Arthur, 1990,

¹⁰Hull (1988) and Campbell (1987) argue that scientific concepts evolve in a manner described by evolutionary models. Their non-technical approach has not been formalized in explicit models, however. Higgs (1987, ch. 4) develops a model with similar dynamics of evolution in his ratchet model of governmental growth.

p. 18). Arthur goes on to point out that this sort of finding is unfamiliar in economics "where our habit of thinking is that if there is a better alternative, it would be chosen" (p.18). A basic problem with this standard substantive rationality result is that the menu of choices is not really known a priori by the chooser. This menu is itself to be learned, and this learning can often involve exploring unknown territory. Such exploration is what Arthur is attempting to model, in a way similar to that proposed by Holland (1975) in his suggestion about a genetic algorithm for the maximization of mathematical functions which standard techniques cannot solve. Arthur speculates at the conclusion that "there is thus an 'ecology' of decision problems in the economy with earlier patterns of decisions affecting subsequent decisions. This interlinkage would tend to carry sub-optimality through from one decision setting to another. ¹¹ The overall economy would then follow a path that is partly decided by chance, is history dependent, and is less than optimal". (p. 19) All that is missing from Arthur's speculation is an explicit recognition of the role of ideologies in this process.

VI. Dynamics of Mental Models and Institutions

Arthur's speculation provides us with a tentative entering wedge to further speculate about the dynamic process of cognitive change occurring as societies and economies evolve. That society's development have been sub-optimal is certainly not open to question. The path-dependence of the institutional development process can be derived from the way cognition and institutions in societies evolve. Both usually evolve incrementally but the latter, institutions, clearly are a reflection of the evolving mental models. Therefore the form of learning that takes place is crucial.

VI.A. Changes in Mental Models as Punctuated Equilibria

The usual modeling of learning in economics involves Bayesian ideas. The Bayesian learner starts out with some sort of prior distribution of beliefs distributed over some pre-defined model space involving the learner's current ideas about how to think about the phenomenon that is the object of the learning. The prior beliefs are updated by some direct learning which generates observational data. This transition of prior beliefs into posterior beliefs, with an unchanging model space is usually thought of as a gradual process with the posterior beliefs some sort of compromise between the peak of the prior beliefs and the model judged most likely by the data alone (Leamer, 1978). This approach to learning misses some crucial features of learning that we believe can be captured by the approach we wish to follow. Bayesian learners are never surprised, or forced within the updating process to completely change the dimensions of the model space. Such surprises or drastically revised models can be interpreted as representational redescrptions, and involve trajectories which can be described as punctuated equilibria of the sort analyzed in Denzau and Grossman (1993).¹²

¹¹By suboptimality here, we mean that there were technically feasible alternatives, implementable in or as humanly feasible institutions or organizations, that would have resulted in higher ex ante and ex post rates of economic growth, without reducing consumption levels.

¹²This approach works at the level of the individual chooser. But many of the changes we wish to understand are social, such as changes in informal institutions or ideologies. We believe that two approaches to this aggregation problem are likely to bear fruit. First, Kuhn (1970) argues that the choice of a new paradigm is at the individual level. In a crisis, individuals choose facing a confusion of evidence and alternative explanations. The resolution of a scientific crisis is an intersubjective decision in the shared mental models of the members of a scientific community come to a consensus on the new basis for their future studies. Kuhn's ideas have been important in our coming to the ideas we now espouse, and there is more to be mined in this approach.

A second approach to the aggregation question involves the recent work of Bikhchandani, Hirshleifer and Welch (1992) on informational cascades. In their model, only a small number of individuals make choices on the basis of their own mental models. The others in the society follow the choices of these decision leaders, free riding

Punctuated equilibrium involves long periods of slow, gradual change punctuated by relatively short periods of dramatic changes, which we can presume to be periods of representational redescription. This reconceptualization is illustrated in two of the graphs in Hutchins and Hazlehurst. These graphs (Figures 7, 10 and 11 in the original) show the results of learning patterns directly, and with the help of cultural artifacts. The cultural learning approximately halves the time required to learn the relation between moon phases and tides. Both direct learning and the culturally mediated learning show patterns of punctuation. For an extended period, neither type of learning enables the pattern to be acquired. Then the probability of successful acquisition through mediated learning starts increasing steeply, as shown in Figure 4, up to about 60%, and slowly increases thereafter. The same pattern, with a less steep slope, is shown for direct learning. Figure 5, reproducing Figure 10 of Hutchins and Hazlehurst (1992), shows the learning results more directly. The mean square error of the learners starts out on a plateau at 0.25, and drops precipitously in less than 10 generations to a new lower plateau at which it remains. Looking for punctuated equilibrium dynamics and the accompanying representational redescription requires research from the viewpoint we present here, and is yet to be performed.

The punctuated equilibrium approach to the dynamics of mental models has implications like those discussed by Kuhn. They differ somewhat because of the crucial difference created by the attempt in science to maintain the precision of terms as opposed to their plasticity in a popularly held and communicated mental model. In science, Kuhn argues that the relatively precise nature of concepts helps keep a paradigm or conceptual framework almost fixed for long periods. But this precision of terms must withstand the evolution of the meaning of terms which continually occurs in popular spheres. Consider the meaning of the Declaration of Independence with its phrase stating the "All men are created equal." The precise meaning of this phrase in Thomas Jefferson's mind when he wrote it is vastly different from the interpretation given it by the Abolitionists 50 years later, or by most Americans today. We leave discussion of this evolution purposely tacit to allow the reader to step through the evolutionary process.

A crucial feature of this sort of evolution is the bringing of new meanings from related mental models by analogy or metaphor. This process is a natural feature of the way our brains generalize and utilize concepts, and intellectual historians such as Carl Lotus Becker (1932, ch.1) have already used the idea of an evolving "climate of opinion" to analyze the changing meaning of terminology and ideological constructs. These gradual and continual changes, in which new meanings in one field of application gradually transfer into another set of mental models, generates the ideological counterpart of Kuhn's normal science. Normal ideology may attempt to resist change, through having ideological scholars and purists, but we expect that ideologies gradually change due to the changing meanings of its terms and concepts in other models, as well as changing use in common parlance. New concepts that have become important parts of the climate of opinion, both to the intellectuals and to the population en masse, can also get brought into the set of ideas in an ideology, as the gradual accommodation of Darwinism suggests.

Let's think further about this process of accommodation and change in shared mental models. The process does not always progress smoothly or easily. Instead, ideological purists, like religious fundamentalists, try to resist any change, and their resistance might generate a crisis. Such a slowing of gradual change through attempts to maintain purity would create an increasing gap between the general climate of opinion

on their efforts. This approach seems to have substantial value in discussing the shared mental models which many people acquire about religion - most people acquire their models by learning from the original texts or from the learned teachers of the doctrine. Changes in the interpretations of these teachers can be acquired as indirect learning through training or just from going along with their interpretations until one acquires the changes in one's own mental models.

and the "pure" ideology. An example of such a gap seems to be occurring in Castro's Cuba today, and is argued by Przeworski (1991, pp. 1-9). When the ideology finally changes, if it does, it would generate a punctuation, i.e., a short, relatively rapid change. Another alternative dynamics would be that gradual evolutionary change and the incorporation of new elements can endogenously generate a crisis for a different reason. The basis for this crisis would be the discovery of a lack of logical consistency in the ideology, or the discovery of a new set of implications which are viewed as disturbing by adherents of the ideology. The communication of this sort of problem could then be used by an entrepreneur to make a punctuated change in the ideology or religion to further the entrepreneur's own goals.

The existence of discovery and surprise is related to another cognitive problem. We simply do not have the abundance of cognitive resources such that our mental models can actually be logically coherent, or be certain that our ideological beliefs are logically coherent.¹³ Many individuals can understand the inconsistency among 3 statements, such as in the Paradox of Evil; i.e., God is desirous of humans living in happiness, God is omnipotent, and evil exists and makes humans unhappy. These three statements, interpreted naturally, result in a logical inconsistency that has been termed the Paradox of Evil which some human religions attempt to resolve. But in the move to 4 statements, as in Arrow's Theorem, we seem to pass across a complexity boundary. Most people find the Arrow result paradoxical, even after being presented with it. That it was only "discovered" in the 1940s suggests its complexity, even though the underlying axioms had been used for some time. The fact that Arrow himself got the proof wrong in both the 1951 and the 1962 editions of the book is suggestive of the complexity of the logical incoherence problem.

If we cannot immediately see, or even at times understand an argument about, the logical incoherence among 4 (or 8) statements, then it is quite likely that logical incoherence could exist in any modestly complex ideology. Demanding such coherence of an ideology, as Hinich and Munger (1992, ch. 1) do, is to talk about a world in which cognitive resources are truly abundant, and not finite, a view that Cherniak terms, "We have God's mind." An ideological entrepreneur who learns of an incoherence or a disturbing implication of the ideology could utilize this in order to help reinterpret that ideology in ways more suitable to the entrepreneur's goals. We believe that some of the numerous religious controversies that have helped create new sects also result from both disturbing applications and incoherence problems.

VII. Summary

We began this essay by noting that it is impossible to make sense out of the diverse performance of economies and polities if one confines one's behavioral assumptions to that of substantive rationality in which agents know what is in their self-interest and act

¹³As Cherniak (1986, p. 79-80) notes, this question of logical consistency is a very difficult problem as characterized by mathematical complexity theory (technically termed an NP-complete problem, not solvable in polynomial time), and the only known algorithms for implementation are exponential time ones. By a calculation of the sort used on any combinatorial explosion problem (like the British Museum algorithm that would allow monkey typists to write all of Shakespeare), it has been shown that it is impossible for all the possible calculation resources of the entire universe computing for all of the time the universe has existed to determine the logical consistency of more than 138 propositions, assuming they all are well-defined. Even if 138 propositions is sufficient for stating an ideology, the cognitive limitations of a finite human agent are vastly more limited than the calculation suggests. Higgs (1987, p. 37) suggests that ideology is "a somewhat coherent, rather comprehensive belief system about social relations." This seems a definition that is much more defensible than requiring logical coherence initially, and at all times, a task that would seem impossible given the plasticity of language.

accordingly. But once we open up the black box of "rationality," we encounter the complex and still very incomplete world of cognitive science. This essay is a preliminary exploration of some of the implications of the way by which humans attempt to order and structure their environment and communicate with each other. Does the argument have relevance for social science theory? Certainly it does. Ideas matter and the way by which ideas evolve and are communicated is the key to developing useful theory which will expand our understanding of the performance of societies both at a moment of time and over time. At a moment of time, the argument implies that institutions and the belief structure are critical constraints on those making choices and are, therefore, an essential ingredient of model building.

Over time, the approach has fundamental implications for understanding economic change. The performance of economies is a consequence of the incentive structures put into place; that is, the institutional framework of the polity and economy. These are in turn a function of the shared mental models and ideologies of the actors. Whether we pursue the framework suggested by Arthur (1992) or the notion of punctuated equilibrium for the dynamics of mental models, we get some common results. The presence of learning creates path-dependence in ideas, ideologies and then in institutions. Arthur argues that a concept discovered by an individual that is useful in explaining the world is more likely to persist in one's mental models, and this implies path-dependence. The same path-dependence is implied by our related evolutionary interpretation. In both approaches, systems of mental models exhibit path-dependence such that history matters, and in both suboptimal performance can persist for substantial periods of time.

References

- Alchian, Armen A., "Uncertainty, Evolution and Economic Theory," *Journal of Political Economy*, , 58 (1950), 211-21.
- Arthur, W. Brian, "Self-reinforcing mechanisms in economics," in P.W. Anderson, David Pines and Kenneth Arrow (eds.), *The Economy as an Evolving Complex System*, Boston: Addison-Wesley Publ., 1988, 33-48.
- _____, "A learning algorithm that mimics human learning," 90-026, Santa Fe Institute Economics Research Program, Nov., 1990.
- _____, "On learning and adaptation in the economy," Institute for Economic Research Discussion Paper #854, Queen's University, May 25, 1992.
- Aumann, R.J., "Survey of repeated games", in *Essays in Game Theory and Mathematical Economics in Honor of Oskar Morgenstern*, ed. R.J. Aumann et al., Zurich: Bibliographisches Institut, p.11-42, 1981.
- Becker, Carl L. *The Heavenly City of the Eighteenth Century Philosophers*, New Haven: Yale University Press, 1932.
- Bikhchandani, Sushil, David Hirshleifer and Ivo Welch, "A theory of fads, fashion, custom, and cultural change as informational cascades," *Journal of Political Economy*, 100(5) (Sept., 1992), 992-1026.
- Binmore, K. "Modeling rational players: Parts I and II" *Economics and Philosophy*, 3, 179-214; 4, 9-55, 1987-88.
- Burns, Penny, "Experience and decision making: A comparison of students and businessmen in a simulated progressive auction," *Research in Experimental Economics*, 3 (1985), 139-57.
- Campbell, Donald T., "Evolutionary Epistemology," by Donald T. Campbell, Raditzky, Gerard and W. W. Bartley, III (eds.), *Evolutionary Epistemology, Rationality, and the Sociology of Knowledge*, La Salle, IL: Open Court, 1987, 47-90.
- Chandler, A.D., Jr., *The Visible Hand: The Managerial Revolution in American Business*, Cambridge: Harvard Univ. Press, 1977.
- Cherniak, Christopher, *Minimal Rationality*, Cambridge: M.I.T. Press, 1986.
- Churchland, Paul M., *A Neurocomputational Perspective: The Nature of Mind and the Structure of Science*, Cambridge: M.I.T. Press, 1989.
- Clark, Andy, *Microcognition: Philosophy, Cognitive Science and Parallel Distributed Processing*, Cambridge: M.I.T. Press, 1989.
- _____, and Annette Karmiloff-Smith, "The cognizer's innards: A psychological and philosophical perspective on the development of thought," *Mind and Language*, (forthcoming).
- Coursey, Don L., and Edward A. Dyl, "Price Limits, Trading Suspensions, and the Adjustment of Prices to New Information," unpubl. ms., Business School, Washington Univ., Feb. 1990.
- Coursey, Don L., and Charles F. Mason, "Investigations concerning the dynamics of consumer behavior in uncertain environments," *Economic Inquiry*, 25 (Oct. 1987), 549-64.
- Denzau, Arthur and Peter Grossman, "Punctuated equilibria: A model and application of evolutionary economic change." unpublished ms., Economics Department, Washington University, April 1993.
- Downs, Anthony, *An Economic Theory of Democracy*, New York: Harper and Row, 1957.
- Easley, David, and John Ledyard, "Theories of price formation and exchange in double oral auctions," in Friedman and Rust (1992).
- Eldredge, N., and S.J. Gould, "Punctuated equilibria: an alternative to phyletic gradualism," in T.J.M. Schopf (ed.), *Models in Paleobiology*, San Francisco: Freeman, Cooper and Co., 1972, 82-115.
- Feldman, Julian, "An analysis of predictive behavior in a two-choice situation," unpubl. Ph.D. dissert., Carnegie Institute of Technology, 1959.

_____, Paul Milgrom, John Roberts, and Robert Wilson, "Rational cooperation in the finitely repeated prisoners' dilemma," *Journal of Economic Theory*, 27 (1982), 245-52.

Kuhn, Thomas S., *The Structure of Scientific Revolutions*, Chicago: Univ. of Chicago Press, 2nd ed., 1970.

LaCroix, Sumner J., "Homogenous middleman groups: What Determines the Homogeneity?" *Journal of Law, Economics, and Organization*, 5:1 (1989) 211-22.

Leamer, Edward E. *Specification Searches: Ad Hoc Inference with Nonexperimental Data*, New York: Wiley, 1987.

Lindgren, Kristian, "Evolutionary phenomena in a simple dynamics," in C. Langton, C. Taylor, J.D. Farmer and S. Rasmussen (eds.), *Artificial Life II*, (Redwood City, CA: Addison-Wesley, 1992), 295-312.

Marks, Robert, "Repeated games and finite automata" in *Recent Developments in Game Theory*, eds. John Creedy, Jeff Borland, and Jurgen Eichberger, Elgar: Brookfield Vt. 1992.

McCaleb, Thomas S., and Richard E. Wagner, "The Experimental Search for Free Riders: Some Reflections and Observations," *Public Choice*, 47 (1985), 479-90.

North, Douglass C., *Institutions, Institutional Change and Economic Performance*, Cambridge: Cambridge Univ. Press, 1990.

Przeworski, Adam, *Democracy and the Market: Political and Economic Reforms in Eastern Europe and Latin America*, Cambridge: Cambridge Univ. Press, 1991.

Sowell, Thomas, *Knowledge and Decisions*, New York: Basic, 1980.

Smith, Vernon L., and James M. Walker, "Monetary rewards and decision cost in experimental economics," *Economic Science Laboratory*, Univ. of Arizona, Nov., 1990.