

# Are there developmental milestones in scientific reasoning?<sup>1</sup>

Anne L. Fay   David Klahr   Kevin Dunbar  
Carnegie Mellon                      McGill  
University                                      University

## Abstract

*This paper presents a conceptual framework that integrates studies on scientific reasoning that have been conducted with different age subjects and across different experimental tasks. Traditionally, different aspects of scientific reasoning have been emphasized in studies with different aged subjects, and the different literatures are somewhat unconnected. However, this separation leads to a disjointed view of the development of scientific reasoning, and it leaves unexplained certain adult behaviors in very difficult scientific reasoning contexts. In this paper we attempt to integrate these three approaches into a single framework that describes the process of scientific reasoning as a search in an hypothesis space and an experiment space. We will present the results from a variety of studies conducted with preschool, elementary school, and adult subjects, and will show how differences in performance can be viewed as differences in the knowledge and strategies used to search the two spaces. Finally, we will present evidence showing that, in sufficiently challenging situations, adults exhibit deficits of the same sort that young children exhibit, even though one might have expected that these developmental milestones were long since passed.*

Experimental studies of the development of scientific reasoning skills have produced three distinct and somewhat disjoint literatures. Studies focusing on what Klayman and Ha (1987) call "positive test bias" (the tendency to seek

<sup>1</sup> Address correspondence to Anne L. Fay Department of Psychology, Carnegie-Mellon University, Pittsburgh, PA 15213, USA. E-mail address: fay@psy.cmu.edu. The first author was supported by a Post-doctoral Fellowship from the James S. McDonnell Foundation Program in Cognitive Studies for Educational Practice. The second author was supported in part by the Personnel and Training Research Program, Psychological Sciences Division, Office of Naval Research, Contract N00014-86K-0349, and in part by grants from NICHD (R01-HD25211-01A1) and the A.W. Mellon Foundation. The third author was supported in part by a grant from the Natural Sciences and Engineering Research Council of Canada, grant number OGP0037356

instances that are expected to confirm one's current hypothesis) have concentrated on adult performance; studies on subjects' faulty strategies for the "coordination of theory and evidence" (Kuhn, 1989) have been conducted primarily with adolescents; and studies examining the understanding of necessity and possibility have been conducted with preschoolers. Rarely is one of these phenomena studied in a different age group (i.e., we know of no studies focusing on adults' understanding of the logic of indeterminacy, nor any of preschoolers' positive test bias.)

One possible justification for the different foci is that there might be a sequence of developmental milestones in the acquisition of a complete set of scientific reasoning skills. If so, then it would be prudent for investigators interested in different age levels to address the most obvious inadequacies of their subjects. However, this separation leads to a disjointed view of the development of scientific reasoning, and it leaves unexplained certain adult behaviors in very difficult scientific reasoning contexts. In this paper we attempt to integrate these three approaches into a single framework that describes the process of scientific reasoning as a search in an hypothesis space and an experiment space. We will present the results from a variety of studies conducted with preschool, elementary school, and adult subjects, and will show how differences in performance can be viewed as differences in the knowledge and strategies used to search the two spaces. Finally, we will present evidence showing that, in sufficiently challenging situations, adults exhibit deficits of the same sort that young children exhibit, even though one might have expected that these developmental milestones were long since passed.

## Components of Scientific Reasoning

Klahr & Dunbar (1988) have conceptualized the process of scientific reasoning as a dual search in an experiment space and an hypothesis space. Figure 1 depicts the two spaces and the logical relations between them. The upper box is the hypothesis space, which consists of specific hypotheses related to the domain. The lower box is the experiment space. Within this space are the experiments that can be conducted in the domain. The arrows connecting the boxes in the two spaces specify how the experimental outcomes bear on the hypotheses. The heavy arrow between Hypothesis A and Experiment 1 indicates that only Hypothesis A is consistent with the outcome of Experiment 1. This reflects a

*Determinate* relation. The light arrows between Hypothesis A and Experiment 2 and between Hypothesis B and Experiment 2 indicate that both Hypothesis A and Hypothesis B are consistent with the outcome of Experiment 2. This reflects an *Indeterminate* relation, whereby the outcome of Experiment 2 cannot discriminate between Hypothesis A and Hypothesis B. The absence of arrows between Hypothesis B and Experiment 1, Hypothesis B and Experiment 3, and Hypothesis A and Experiment 3, indicates that these hypotheses are inconsistent with each of these outcomes, reflecting an *Impossible* relation.

The goal of scientific discovery is to generate experiments and hypotheses that will eventually result in a relation of determinacy, whereby only one hypothesis remains consistent with all the experimental outcomes.<sup>2</sup> Thus, the components of scientific reasoning consist of: 1) Identification and understanding of the relations between experiments and hypotheses (i.e., understanding the logic of necessity and possibility); 2) Generation of informative experiments, (i.e., generating experiments that further specify the relations between the two spaces so as to prune the search in the space of hypotheses); 3) Hypothesis generation and revision (i.e., generation of hypotheses from either analogy or via induction from experimental outcomes.)

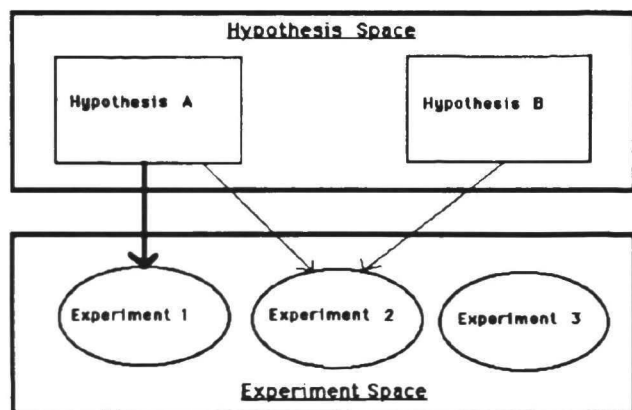


Figure 1. Schematic representation of the two spaces and the relations between them. Heavy arrow indicates a determinant relation between experimental outcome and hypothesis. Light arrows indicate an indeterminate relation and no line indicates an impossible relations.

These three general abilities are fundamental to

<sup>2</sup> In real-world situations, one never has a determinate relation, as there are an infinite number of possible hypotheses, and there is always the possibility that new evidence will disconfirm the current hypothesis. Nonetheless, the goal of science can be seen as the elimination of all *current* competing hypotheses until only one remains consistent with the existing evidence.

the process of scientific reasoning, and their deficits are characteristically associated with specific age groups. Adults recognize and understand the implications of indeterminacy, and have heuristics for designing informative experiments, but are notoriously biased toward confirmation in rule discovery tasks (Gorman & Gorman, 1984; Gorman, Stafford & Gorman, 1987; Wason, 1960, 1968).<sup>3</sup> Adolescents, like adults, understand the notion of indeterminacy, but in addition to their bias toward confirmation, they lack the strategies and knowledge for designing informative experiments. Preschool children demonstrate all these deficits, but also show a failure to recognize and/or understand the implications of determinate vs indeterminate situations. Thus, acquisition of these three abilities might be viewed as *milestones* in the development of scientific reasoning skills. However, the picture is not that straightforward, as we shall argue below. First, however, we will further elaborate each of the three components enumerated above.

#### Identifying the relation between experiments and hypotheses

One of the basic components of scientific reasoning is the ability to recognize and understand the implications of confirming, and disconfirming evidence. Understanding the implications of these conditions is based on the distinction between determinate and indeterminate outcomes. In Indeterminate situations, the evidence is insufficient to discriminate one hypothesis from another. Until this concept is available, the process of scientific discovery will be severely flawed. Failing to recognize a situation as indeterminate will result in premature termination of the generate-experiment process because a confirming instance will be erroneously identified as sufficient to accept a theory.

Research with preschool children has shown they lack the concept of indeterminacy. In terms of Figure 1, they fail to realize the relation between Experiment 2 and the Hypothesis space. In a study extending Pieraut-Le Bonnic's (1980) investigations of children's understanding of possibility and necessity, Fay & Klahr (1990), presented kindergarten children with two boxes of building materials, and a series of objects, one at a time, made from materials taken entirely from one box or the other. For example, Box A might contain sticks and curves and Box B might have sticks and squares. A probe object comprised of

<sup>3</sup> But see Farris & Revlin (1989) for a novel reinterpretation.

only sticks would be *indeterminate* because it could have been constructed from either box. A *determinate* probe object would be one constructed from sticks and curves (only Box A could have been used to make it). The children were asked whether they could tell which box was used to make the object .

Figure 2 shows a schematic of the problem. As can be seen, the task can be mapped directly onto the components shown in Figure 1, with the boxes representing hypotheses and the probe objects representing experiments. Thus in this context the child is presented with a finite hypothesis space (e.g. box with sticks and curves vs. box with sticks and squares) and an experimental outcome (stick & curve object vs. stick-only object) but must determine the relations that exist between them (stick & curve object is determinate vs. sticks-only is indeterminate). All the children correctly identified the determinate situation, but only 53% consistently identified the indeterminate situation. Those who failed to recognize the indeterminate situation misidentified it as determinate, claiming that they *could* tell for sure which box had been used to construct the indeterminate probe object.

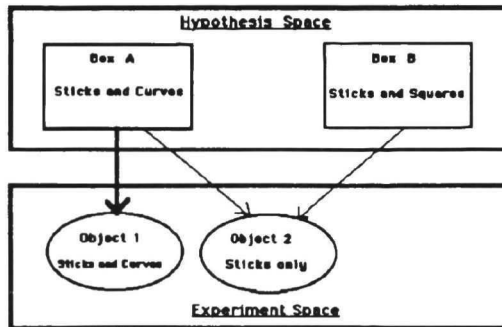


Figure 2. Experimental setup for the Possibility-Necessity study. Heavy arrows indicate a determinate relation and lights arrows indicate an indeterminate relations.

This failure can be explained, in part, by a lack of understanding of the concept of logical necessity. Evidence for this interpretation comes from children's justifications for their responses on determinate problems, which were coded as being based on either positive or negative reasoning. Children were scored as using positive reasoning if their justification was based on the confirmatory relation between the determinate box and the object (e.g. "you used this box because it has sticks"). They were scored as using negative reasoning if their justification referred to the necessity of disconfirmation of the other response (e.g."you had to use this box

because the other box doesn't have any curves"). Negative reasoning implies an understanding of logical necessity. That is, it suggests that the child recognizes the insufficiency of confirmation alone and therefore searches the entire hypothesis space to determine the other experiment-hypothesis relationships. Table 1 shows the relation between the type of reasoning used on the determinate problems and performance on the indeterminate problems. The results suggest that the tendency to use negative reasoning is related to the recognition of indeterminate situations.

TABLE 1: Children who use negative reasoning on determinate problems are more likely to be correct on indeterminate problems.

PERFORMANCE ON INDETERMINATE PROBLEMS  
(cell entries are number of responses)

REASONING ON DETERMINATE PROBLEMS	Incorrect	Correct
Positive Reasoning	22	20
Negative Reasoning	6	19

(Chi Square=5.19, p<.025)

Young children demonstrate a lack of understanding of logical necessity, a prerequisite of scientific reasoning. Failing to recognize an indeterminate situation, or to understand its implication, will result in a premature termination of the search based on finding a confirmatory relation between data and theory.

Generating Informative Experiments

The ability to recognize the relations between hypotheses and experiments can be seen as a prerequisite for the skill of generating informative experiments. Informative experiments are designed for the purpose of pruning the hypothesis tree, that is, eliminating impossible hypotheses, and reducing the set of consistent hypotheses. In this situation, the subject is provided with an hypothesis, or enters with a prior hypothesis, and must generate experiments which will lead to the confirmation or disconfirmation of the hypothesis. In reference to Figure 1, the subject is provided with an hypothesis (e.g. Hypothesis B), and the task is to generate experiments that will either disconfirm the hypothesis (e.g. E3), or will discriminate between existing confirming hypotheses (e.g. Experiment 1). Subjects want to avoid writing experiments that fail to discriminate between existing hypotheses (e.g. Experiment 2), or at the least, recognize them as being indiscriminating. Thus, this ability is dependent upon the ability to recognize and understand the relations between

experiments and hypotheses. In addition, it involves an understanding of the goal of experiment generation (reducing the hypothesis tree), and the skills for constructing experiments that will serve these goals.

In a series of experiments with children (8 to 13 years old), and adults, we examined subjects' ability to generate informative experiments (Klahr, Dunbar & Fay, 1990; Klahr, Fay & Dunbar 1990). Subjects were trained to operate a simple programmable device by entering commands (for moving forward, backward, turning right and left, and firing its cannon) and then pressing a GO key to execute the program. This would move an icon on a workstation screen according to the program the subject had entered.<sup>4</sup> Once trained to criterion, they were then asked to discover how an additional function, the REPEAT key, worked. They were then provided with an hypothesis (which was always incorrect), and were asked to write programs to find out if the hypothesis was correct or, if it wasn't, to find out how REPEAT worked.

The design of the study (See Table 2) crossed the plausibility of the given and actual hypotheses. Thus, subjects could be given either a highly plausible or highly implausible hypotheses for how REPEAT worked, and the device was actually programmed to interpret REPEAT in some different, but either plausible or implausible, way. Subjects were given one of the rules and the device actually worked according to a different rule. This effect of these given-actual hypotheses conditions will be expanded on in the following section.

**Table 2: Design of "negative feedback" study**

	Actual	
	Plausible	Implausible
Given Plausible	Theory refinement	Theory replacement
Implausible	Theory replacement	Theory refinement

Overall, children performed poorly in discovering the correct rule. Only one-third of the younger children, and half of the older children discovered the rule, compared to 83% of the adults.

<sup>4</sup>This "microworld" was a simulated version of the BigTrak, a programmable robot toy that moved around on the ground, originally used by Shrager & Klahr, 1986.

One contributing factor to this trend is the degree of informativeness of the programs that the subjects wrote. First, children appear to differ from adults in terms of their awareness of the goals of experimentation. Whereas 83% of the adults made statements referring to experimental design goals, only 20% of the younger children and 47% of the older children made such comments. The quality of these statements also differed. The adults stated experiment goals in terms of increasing the observability and informativeness. The youngest children, on the other hand, primarily made output goal statements (e.g. move it in a square) and some observability goals (e.g. use N=1 otherwise it's too confusing). The older children focused on observability goals (e.g. shorten programs, use easily traced commands).

The above data is based on verbal reports, and as such the tendency to verbalize may be different for the different age groups. A second analysis examined the types of programs that were written. The experimental space for this problem can be viewed along two dimensions, one dimension being the number of commands in the program ( $\lambda$ ) and the other being the magnitude of the argument for REPEAT (N). The ideal experiment is one which maximizes the informativeness of the outcome while minimizing the complexity (i.e. maximizing observability or interpretability of outcome). In the current setting, this means writing minimum length programs that can discriminate the effect of the REPEAT function. By this criterion, the "best" program has a length ( $\lambda$ ) of 3 and a REPEAT argument value (N) of 2. The three age groups differed in their tendency to write programs with these properties. Compared to a random model, the children were 1.5 times *less* likely to generate the ideal experiment whereas the adults were 5 times *more* likely to run such an experiment. In addition, adults were much more systematic in the way that they moved in the experiment space. Their experiments had more of the flavor of a careful experimental series than did the children's.

In summary, children appear to lack the knowledge required to generate informative experiments. Part of this deficit involves a failure to understand the goals of experimentation. Whereas adults' goals were directed toward informativeness and observability, children's goals were directed toward producing a desired effect, and, for the older children, observability. However, even

though the older children recognized the importance of observability, they were not overly successful at designing interpretable experiments.

### Generating and revising hypotheses

The final component of scientific reasoning is the ability to generate and revise hypotheses in response to experimental outcomes. Combined with the other abilities, this situation can be depicted in Figure 1 by having no boxes specified or present in either the experimental or hypothesis space. Thus all the components of the task must be generated by the subject. In series of studies using the physical BigTrak device, adults and third to sixth grade students were trained on all the functions of the device except the REPEAT and were then asked to write programs to figure out how REPEAT worked (Klahr & Dunbar, 1988; Dunbar & Klahr, 1988). There are two main differences in these studies as compared to the studies mentioned in the previous section: First, in these studies subjects were not given any hypothesis, and had to generate their own hypotheses from the start, and second, there was only one rule for REPEAT, it caused the device to repeat the last N instructions once, where N refers to the argument for REPEAT. There was a strong age effect: 19 of the 20 adults, but only 2 of the 22 children successfully discovered the correct rule, although over half of the children *believed* that they had correctly identified it. Part of the failure can be attributed to the non-informative programs that the children wrote. However, both the adults and the children had the same proportion of experiments from the most informative region of the experiment space, where  $\lambda > N$  and  $N > 1$ . Based on the outcomes from experiments conducted in this region adults were able to induce the correct hypothesis but the children were not.

The hypotheses that were generated by the subjects in this study can be classified into two frames based on the function of N. One frame, Counter-Frame, assigned a role to N where the number indicated how many times something got repeated. The other frame, Selector-Frame, assigned a role to N where the number indicated which instructions got repeated. Children and adults demonstrated an initial preference for the Counter-Frame hypotheses, (which is incorrect in this study). But, whereas adults were able to abandon this hypothesis frame in light of disconfirming evidence and generate the Selector-Frame, the children tended to maintain Counter-

Frame hypotheses in spite of disconfirming evidence. Thus the children's search of the hypothesis space was constrained to Counter hypotheses.

Further evidence for this comes from the negative feedback studies described earlier, in which subjects were given an initial hypothesis that could be either from the same frame as the actual hypothesis or from the other frame as the actual hypothesis. Figure 3 shows the effect for Given-Actual Hypothesis conditions. The children were successful when the device worked as a (plausible) Counter, but failed to get the correct rule when it worked as a (implausible) Selector.

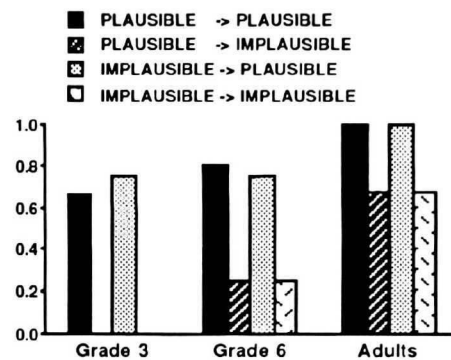


Figure 3. Proportion of subjects discovering correct rule when given a Plausible (counter) or implausible (selector) hypothesis and actual rule was Plausible (counter) or Implausible (Selector).

Prior to running any programs, children and adults differed in terms of their willingness to entertain a Selector hypothesis. Table 3 shows the proportion of subjects in each age group that initially accepted the Given hypothesis or one from the same frame as the Given. The children, especially the younger ones, find the Selector hypothesis very implausible. More than half of the children who rejected the Selector proposed a Counter hypothesis instead. Adults, on the other hand, demonstrated some skepticism over their Given Selector hypothesis, but rather than reject it, they proposed other hypotheses *in addition to* it, and these alternative hypotheses were most likely to be Counters.

Table 3. Proportion of subjects accepting Given Frame prior to running first experiment.

Group	Given Frame	
	Counter	Selector
3rd Grade	1.00	.12
6th Grade	.89	.63
Adults	1.00	1.00

Children's strategies and goals for searching the hypothesis space appear to be different from adults. Children's prior hypotheses constrain their search of the hypothesis space to those areas they consider plausible, in this case, counter hypotheses. Although adults also have prior hypotheses, their search of the hypothesis space is not so constrained, and they will entertain the possibility of implausible hypotheses. Thus adults will abandon a more plausible hypothesis frame given disconfirming evidence, and search the space for a new, less plausible frame, children continue to search within the plausible frame for particular hypotheses that will explain the experimental outcomes.

### **Milestones or fragile acquisitions?**

Given the characteristic deficits associated with each age range, and given the logical necessity for each of the three skills to be in place before the next one can be reliably assessed, it is tempting to view this as a sequence of developmental milestones, in which a skill, once acquired, can be reliably invoked in a wide range of situations and can provide the basis for the subsequent acquisition. However, in other careful analysis of children's strategy acquisition (e.g, Siegler and Jenkins, 1989), it has been shown that the story is not so simple. A new strategy or skill may appear for a while, and then disappear for a protracted period. Or a strategy that seemed quite robust, may, in contexts of sufficient complexity, be abandoned, as subjects revert to simpler, and inadequate, strategies. In the domain of scientific reasoning, we have found just this situation. Dunbar (1989) found that strong prior beliefs about hypotheses can overly constrain search of the hypothesis space, and produce behavior that, at its core, reveals a severely limited ability to discriminate determinate from indeterminate outcomes.

Adult subjects were given training in a simulated molecular genetics laboratory and were shown how to go about discovering how certain genes control the enzyme production of other genes by switching them on when a nutrient is present. This mechanism was *activation*. Subjects were shown the different variables that could be manipulated (e.g. amount of nutrient present, genetic mutations), and how they could use this information to run experiments and induce the control mechanism. Subjects were then given a new set of genes and were asked to discover how

the enzyme producing genes were controlled. However, the mechanism in this set was *inhibition*: controller genes turn other genes off until a nutrient is present. This can be compared to the Given-Plausible, Actual-Implausible hypothesis condition in the simulated BigTrak studies, as shown in the top-right cell of Table 2. Only 25% of the subjects discovered the inhibition mechanism, similar to the success rate of 6th grade students in the BigTrak study. Subjects often conducted experiments that could have been consistent with many hypotheses, but interpreted the results as confirming their prior (plausible) hypothesis. Sixty-five percent of the subjects remained within the Activation-frame of the hypothesis space, despite experimental evidence that disconfirmed this frame. Thus, like the children in the previous study, their prior hypothesis overly constrained their search of the hypothesis space and also affected their search of the experiment space.

### **Conclusion**

The child-as-scientist view suggests that children go about the world gathering information and building theories (Brewer & Samarapungavan, in press; Karmiloff-Smith, 1988). Other researchers argue that although children may generate theories of their worlds, the process of theory generation and revision is different from that of adults (Kuhn, 1989). The view presented here is that children of different ages have certain characteristic conceptual deficits, which limit their ability to engage in the process of scientific reasoning. We have attempted to show how the three relatively diverse literatures on scientific reasoning can be integrated into a single framework that views discovery as a dual search in a space of hypotheses and experiments, but we have cautioned against a simple view of developmental milestones because of the tendency for people to regress to earlier deficits in sufficiently complex situations. Perhaps this tendency to regress accounts for the substantial educational and institutional supports that provide practicing scientists with the means to maximize the rationality and effectiveness of their efforts at scientific discovery.

## References

- Brewer, W. F. & Samarapungavan, A. (in press). Child theories versus scientific theories: Differences in reasoning or differences in knowledge? In R. R. Hoffman & D. S. Palermo (Eds.), *Cognition and the symbolic processes: Applied and ecological perspectives* (Vol. 3). Hillsdale, NJ: Erlbaum
- Dunbar, K. (1989). Scientific reasoning strategies in a simulated molecular genetics environment. *Proceedings of the 11th annual meeting of the Cognitive Science society*, 426-433. Ann Arbor, MI: Lawrence Erlbaum Associates.
- Dunbar, K., & Klahr, D. (1989). Developmental differences in scientific discovery strategies. In D. Klahr & K. Kotovsky (Eds.), *Complex information processing: The impact of Herbert A. Simon*. Hillsdale, NJ: Erlbaum.
- Farris, H.H., & Revlin, R. (1989). Sensible reasoning in two tasks: Rule discovery and hypothesis evaluation. *Memory & Cognition*, *17*, 221-232.
- Fay, A. L. & Klahr, D. (1990). Cognitive precursors to scientific reasoning: the development of the concepts of possibility and necessity. Working Paper, Dept. of Psychology, Carnegie-Mellon University.
- Gorman, M.E., & Gorman, M.E. (1984). A comparison of disconfirmatory, confirmatory, and control strategies on Wason's 2-4-6 task. *Quarterly Journal of Experimental Psychology*, *36a*, 629-648.
- Gorman, M.E., Stafford, A., & Gorman, M.E. (1987). Disconfirmation and dual hypotheses on a more difficult version of Wason's 2-4-6 task. *Quarterly Journal of Experimental Psychology*, *39a*, 1-28.
- Karmiloff-Smith, A. (1988) The child is a theoretician, not an inductivist. *Mind and Language*, *3*, 183-195.
- Klahr, D., & Dunbar, K. (1988). Dual space search during scientific reasoning. *Cognitive Science*. *12*, 1-48.
- Klahr, D., Dunbar, K. & Fay, A. L. (1990) Designing good experiments to test "bad" hypotheses. In J. Shrager & P. Langley (Eds.), *Computational models of discovery and theory formation*. Morgan-Kaufman.
- Klahr, D., Fay, A.L., & Dunbar, K. (1990). *Developmental differences in experimental heuristics*. Working paper, Department of Psychology, Carnegie Mellon University..
- Kuhn, D. (1989) Children and adults as intuitive scientists. *Psychological Review*, *96*, 674-689.
- Klayman, J., & Ha, Y. (1987). Confirmation, disconfirmation and information in hypothesis testing. *Psychological Review*, *94*, 211-228.
- Pieraut-Le Bonniec, G. (1980). *The Development of Modal Reasoning*. New York, NY: Academic Press.
- Shrager, J., & Klahr, D. (1986). Instructionless learning about a complex device. *Journal of Man-Machine Studies*, *25*, 153-189.
- Siegler, R. S. & Jenkins, E. (1989) *How children discover new strategies*. Hillsdale, N. J.
- Wason, P. C. (1960). On the failure to eliminate hypotheses in a conceptual task. *Quarterly Journal of Experimental Psychology*, *12*, 129-140.
- Wason, P.C. (1968). On the failure to eliminate hypotheses: A second look. In P.C. Wason & P.N. Johnson-Laird (Eds.), *Thinking and Reasoning*. Middlesex: Penguin Books.