# Flawed Foundations of Associationism?

## Comment on Machado and Silva (2007)

C. R. Gallistel
*Rutgers—The State University of New Jersey*

*A. Machado and F. J. Silva (2007) have spotted an important conceptual problem in scalar expectancy theory's account of the 2-standard-interval time-left experiment. C. R. Gallistel and J. Gibbon (2000) were aware of it but did not discuss it for historical and sociological reasons, owned up to in this article. A problem of broader significance for psychology, cognitive science, neuroscience, and the philosophy of mind concerns the closely related concepts of a trial and of temporal pairing, which are foundational in associative theories of learning and memory. Association formation is assumed to depend on the temporal pairing of the to-be-associated events. In modeling it, theorists have assumed continuous time to be decomposable into trials. But life is not composed of trials, and attempts to specify the conditions under which two events may be regarded as temporally paired have never succeeded. Thus, associative theories of learning and memory are built on conceptual sand. Undeterred, neuroscientists have defined the neurobiology-of-memory problem as the problem of determining the cellular and molecular mechanism of association formation, and connectionist modelers have made it a cornerstone of their efforts. More conceptual analysis is indeed needed.*

*Keywords:* memory, associative learning, learning distributions, temporal pairing, learning trials

I am deeply sympathetic to the call for more sustained and penetrating conceptual analysis. The critical conceptual analysis of scalar expectancy theory's (SET's) explanation of the two-standard-interval time-left paradigm that Machado and Silva (2007) described is on target. John Gibbon and I, while working on a long theoretical article (Gallistel & Gibbon, 2000), discussed this problem. None of our discussion found its way into our review, however, for sociological reasons: Behaviorist perspectives dominated the study of basic learning processes for many years. Although behaviorism has been in retreat for decades, the study of animal learning is something of a refuge. Our models imputed to pigeons, rats, rabbits, and mice representational resources and computational capacities that an audience with lingering behaviorist sympathies was likely to find a priori implausible. As far as Gibbon and I could see, an attempt to address the problem would require positing that the animals (a) compared the observed distribu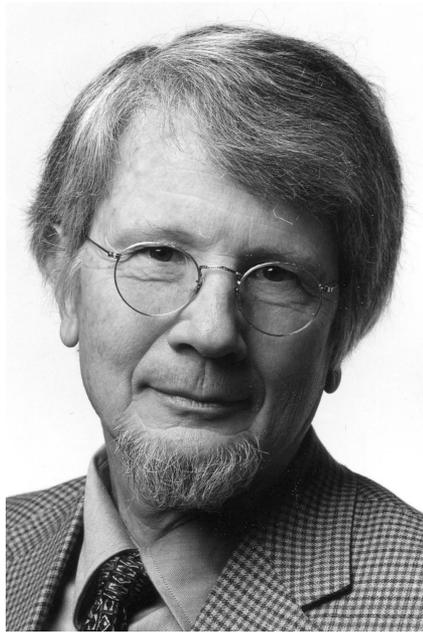tion of intervals with the distributions expected on the random-rate (no temporal structure) hypothesis and the (single) fixed-interval hypothesis (a monotonic normal distribution), (b) derived from this comparison the value of a decision variable, and (c) used that value together with a (presumably) innate function specifying its expected distribution to decide what kind of distribution the intervals they had experienced came from (exponential, unimodal, or bimodal). Putting into print this suggestion about what goes on in the brain of a pigeon, in a source likely to be widely read, would have called our sanity into question, so we did not do it.

In a book based on this same collaboration (Gallistel & Gibbon, 2002), which we correctly assumed few would read, Gibbon and I did call attention to just this suggestion. There, we offered the model just sketched. If the realization that there is a substantial conceptual problem posed by the fact that pigeons recognize a bimodal distribution when they encounter it ever becomes widespread, then we can claim credit for having been among the first to recognize the problem and to suggest an explanation. Meanwhile, our suggestion reposes in an unread book, doing our reputations for sanity no harm.

A behaviorist would object that the model just sketched is vastly more complicated than is required to explain the phenomenon (a bimodal response to a bimodal distribution). The behaviorist would posit an array of internal stimulus traces (neural activities triggered by the conditioned stimulus), each rising and subsiding with a different time course (Grossberg & Schmajuk, 1991; Miall, 1996). When reinforcement occurs after 15 s (the first mode), an association forms between the response and the traces that peak near 15 s. When reinforcement occurs after 240 s (the second mode), an association forms between the reinforced response and traces that peak near that much longer interval.

A problem for the behaviorist account is that the experiment in question (Brunner, Gibbon & Fairhurst, 1994) shows that the pigeons extract three different temporal decision criteria from the bimodal distribution they experience when they end up on the "standard" key, where they encounter one of two delays: 15 s or 240 s. During the

**C. R. Gallistel**

initial part of each trial in this complex experiment, the pigeon was free to switch between two response keys. Sooner or later, at the unpredictable moment of commitment, the key not being pecked at that moment went dead, forcing the bird to finish out the trial on the key it was pecking at the moment of commitment. The first key was called the "standard" key; the second key was called the "time-left" key. The delay from the moment of commitment to reward that the pigeon faced if it ended up committed to the time-left key was the difference (the time left) between an initially very long (hence, unattractive) delay and the time elapsed up to the moment of commitment. The longer the trial went on prior to the moment of commitment, the shorter the time left; hence the more attractive this key became. If the pigeon ended up committed to the standard key, it faced a delay of 15 s on a random half of such trials and a delay of 240 s on the other half. The attractiveness of this key did not vary with time. The datum of principal interest was the time left at which the pigeon switched from the initially more attractive standard key to the increasingly attractive time-left key. That time indicated what the pigeon took to be the expectation of a distribution composed in equal parts of 15-s and 240-s delays. On the trials when the pigeon ended up committed to the fixed-delays key, its pecking probability peaked at 15 s and then (if there was no reward) subsided, only to rise again at 240 s. On trials where it had switched to the time-left key before the moment of commitment, the average time when it switched was at the harmonic mean of these two delays $2/[(1/15)+(1/240)] = 28$. Why it takes the expectation to be the harmonic mean rather than the arithmetic mean is a mystery, but one with ample empirical support, because this is an instance of the ubiquitous hyperbolic discounting. The conceptual challenge is to explain all three facts (the peaks in pecking probability at 15 s and 240 s and the switches at the harmonic mean of these two delays). It is not a small challenge.

I give this embarrassing account of why we chose not to confront a conceptual problem we knew was there, because it illustrates one cause for scientists' reluctance to engage in more conceptual analysis: Careful conceptual analysis tends to show that the explanatory power of our favorite models is an illusion. And it tends to force us to consider assumptions or assertions that we are loath to make, first because their novelty makes them harder to work with and, second, because they almost always seem more complex than the assumptions we have grown comfortable with. That is why scientists tend to stick to their traditional explanations even when they have quite fundamental conceptual flaws that intrude to varying degrees into awareness. It also explains why it is hard to make a scientific living pointing out those conceptual problems.

I write this with some feeling because I have been trying without success for some years to get students of learning and memory to focus on what I take to be a serious conceptual flaw in the foundations of associative learning theory. Associative learning theory is ubiquitous in cognitive psychology, cognitive science, and connectionist modeling. It is also the theory that guides the efforts of investigators all over the world to find the neurobiological basis of learning and memory. Thus, there should be quite general concern with its conceptual foundations.

A cornerstone of those foundations is the concept of a discrete trial, during which two stimuli either are or are not paired, which pairing or failure thereof either strengthens or weakens the association between them.[1] The problem is that life is not—at least on the face of it—composed of trials. The pessimist would say it is just one long trial. Setting that bleak assessment aside, the fact remains that unless a principled way can be found to identify the theorists' conception of a trial with the extraordinarily various—and temporally continuous!—reality to which associative theory is supposed to apply, associative theory has no empirical bite. Absent a means of identifying what it is in the flow of events that constitutes a trial, a theory in which the relevant psychological or neurobiological changes are triggered by what happens on successive trials cannot be made to apply.

At least for myself, the conceptual problem posed by the notion of a trial was first brought into clear view in Rescorla and Wagner's (1972) seminal revision of associative theory, one of the most widely cited and influential papers of the last half century. They were concerned in this paper to explain two experimental results that appeared to challenge the foundations of associative learning theory as then understood. One was the result of the Rescorla (1968)

---

[1] The association, that is, the conductive connection, must, of course, form not between the events themselves but rather between representations or tokens of them internal to the mind or brain. Because it is cumbersome and wordy to keep making this distinction between the events and their representation, I follow the nearly universal practice of writing as if the events themselves were associated.

experiment, which distinguished for the first time between temporal pairing and contingency. For several groups of rats, a noise sounded for 2 minutes at unpredictable times, with an average of 10 minutes between soundings. For one of the groups, shocks occurred only when the noise was present; for another, shocks occurred during the noise, just as for the first group, but they occurred just as frequently during the "background" intervals when the noise was absent.

The temporal pairing of noise and shock, which associative theory takes to be what is critical for association formation, was the same for both groups. However, for the first group, the noise was informative (shock occurred only in its presence), whereas for the other group it was uninformative (the shock was equally likely in its absence). The challenging result was that the second group did not develop an association between the noise and the shock (or at least they gave no behavioral evidence of having done so), whereas the first group did. The implication was that it was the information provided by the noise about the shock that drove the associative process, not the temporal pairing between noise and shock. It was taken for granted that the learning process must be an associative one, that is, that it must involve the formation in some sense of a conductive connection. The conceptual problem posed when this result is viewed from an associative perspective was nicely put by Rescorla (1972, p. 10): "We provide the animal with individual events, not correlations or information, and an adequate theory must detail how these events individually affect the animal. That is to say that we need a theory based on individual events."

Rescorla and Wagner's (1972) hugely influential approach to this problem was to assume (a) that the subjects in the second group developed an association between the shock and the experimental chamber and (b) that this association between the shock and the "background" blocked the formation of an association between the shock and the noise by gaining for itself all of a limited amount of total associative strength. Rescorla and Wagner saw no conceptual problem with this approach until they tried to make a computer simulation of the process by which the shock's association with the background crowded out the association between the noise and the shock. At that point, they confronted the problem of defining what constituted the trials that led to the formation of the association between the chamber and the shock.

Experimentalists generally assume that each occurrence of a potentially predictive (informative) stimulus constitutes a trial. They arrange their experiments so that the predicted event occurs at most once during any one occurrence of the predicting stimulus. In Rescorla's (1968) experiment, each "occurrence" of the chamber stimulus— that is, each experimental session—lasted two hours, during which the rat experienced several shocks at unpredictable times. One could naively assume that each session constituted one trial, hence one pairing of the chamber and shock (ignoring the fact that there were many shocks during that one "trial"). That assumption does not yield the result that Rescorla and Wagner (1972) were after. On that

assumption, trials involving noise and shock occur much more often than trials involving the chamber alone and shock, which leads to the noise–shock association crowding out the chamber–shock association.

To make their account go through, the chamber–shock trials had to be more frequent than the noise–shock trials. Rescorla and Wagner (1972) achieved this by what they acknowledged was an ultimately indefensible ad hoc assumption, namely, that when a rat was in the chamber, the rat had an internal trial timer that, mirabile dictu, parsed its experience of the box into a sequence of 2-minute-long trials. Because there were on average 10 minutes when the noise was not present for every 2 minutes when it was, there were five times as many "trials" with only the chamber and the shock (background-alone "trials") as there were "compound trials," during which the noise (together with the chamber) was paired with the shock.

As a theorist, I do not begrudge Rescorla and Wagner (1972) an ad hoc assumption or two, even ones as implausible as this one (which assumes not only an internal trial timer but one designed or evolved with Rescorla's experiment in mind, because the trials it times are 2 minutes long). Every theorist needs ad hoc assumptions. Sometimes they deal with matters that are pretty clearly outside the scope of the theory. However, in the present case, the assumption was made in order to make the notion of a temporal pairing between two events (being in the chamber and getting shocked) applicable to the experiment being analyzed. The notion of a temporal pairing is at the core of associative theory. (See Colwill, Absher, & Roberts's [1988] discussion of the implications of her demonstration of background conditioning in *Apylsia californica* for Kandel & Schwartz's [1982] model of the cellular/molecular mechanism of learning.) Two events are said to be temporally paired just in case they occur on the same trial. Thus, the concept of a trial is a core concept. No trials, no theory.

It is reasonable to make indefensible ad hoc assumptions concerning core concepts as a temporary expedient, to see how far a theory carries one once one has gotten round the conceptual obstacle that has one buffaloed. However, if one's theory is then to be taken as a true story about physical reality—and associative theory is routinely taken as just that by contemporary neurobiologists interested in the cellular and molecular bases of learning—one is obliged to revisit indefensible ad hoc assumptions at its foundations and find defensible alternatives to them. The phenomenon of background conditioning is experimentally well-established. Moreover, the assumption that conditioning to the context occurs is critical to most contemporary associative explanations of many other phenomena (e.g. Bouton, Westbrook, Corcoran, & Maren, 2006; Denniston, Savastano, & Miller, 2001).

The conceptual problem that Rescorla and Wagner (1972) confronted arises whenever background conditioning is seen, which is to say, almost everywhere, and it lurks in less obvious ways in the shadows surrounding many other attempts to bring associative theory to bear in understanding learning. Connectionist modelers are oblivious to the problem because they structure the input to their models

in trials, thereby avoiding it. It has been 35 years since Rescorla and Wagner overcame the problem with a transparently indefensible assumption, an assumption that I have never seen them or anyone else defend in print as anything but a temporary expedient. In that time, we have not learned to solve the problem, only to ignore it.

## REFERENCES

Bouton, M. E., Westbrook, R., Corcoran, K. A., & Maren, S. (2006). Contextual and temporal modulation of extinction: Behavioral and biological mechanisms. *Biological Psychiatry, 60,* 352–360.

Brunner, D., Gibbon, J., & Fairhurst, S. (1994). Choice between fixed and variable delays with different reward amounts. *Journal of Experimental Psychology: Animal Behavior Processes, 20,* 331–346.

Colwill, R. M., Absher, R. A., & Roberts, M. L. (1988). Context–US learning in *Aplysia californica. Journal of Neuroscience, 8,* 4434–4439.

Denniston, J. C., Savastano, H. I., & Miller, R. R. (2001). The extended comparator hypothesis: Learning by contiguity, responding by relative strength. In R. R. Mowrer & S. B. Klein (Eds.), *Handbook of contemporary learning theories* (pp. 65–116). Mahwah, NJ: Erlbaum.

Gallistel, C. R., & Gibbon, J. (2000). Time, rate, and conditioning. *Psychological Review, 107,* 289–344.

Gallistel, C. R., & Gibbon, J. (2002). *The symbolic foundations of conditioned behavior.* Mahwah, NJ: Erlbaum.

Grossberg, S., & Schmajuk, N. A. (1991). Neural dynamics of adaptive timing and temporal discrimination during associative learning. In G. A. Carpenter & S. Grossberg (Eds.), *Pattern recognition by self-organizing neural networks* (pp. 637–674). Cambridge, MA: MIT Press.

Kandel, E. R., & Schwartz, J. H. (1982, October 29). Molecular biology of learning: Modulation of transmitter release. *Science, 218,* 433–443.

Machado, A., & Silva, F. J. (2007). Toward a richer view of the scientific method: The role of conceptual analysis. *American Psychologist, 62,* 671–681.

Miall, C. (1996). Models of neural timing. In M. A. Pastor & J. Artieda (Eds.), *Time, internal clocks and movement* (Advances in Psychology series, Vol. 115, pp. 69–94). Amsterdam: North-Holland/Elsevier Science.

Rescorla, R. A. (1968). Probability of shock in the presence and absence of CS in fear conditioning. *Journal of Comparative and Physiological Psychology, 66,* 1–5.

Rescorla, R. A. (1972). Informational variables in Pavlovian conditioning. In G. H. Bower (Ed.), *The psychology of learning and motivation* (Vol. 6, pp. 1–46). New York: Academic Press.

Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F. Prokasy (Eds.), *Classical conditioning II* (pp. 64–99). New York: Appleton-Century-Crofts.