

【講演1】

Robert A. Rescorla (ペンシルバニア大学名誉教授)

「Measuring changes in associative learning」

Perhaps the most elementary observation in associative learning is that when you change the relationship between two events, observing organisms change their behavior to reflect that change. This is true whether the relationship you're talking about is between two events which are outside in the world, such as they are in Pavlovian conditioning or whether the relationship is between two events one of which is an animal's behavior and the other of which is in the world, as it is instrumental learning. It's true whether the change that you're arranging involves the institution of a new relationship between two things, as it does in acquisition, or the change involves the destruction of an old relationship, as it does in extinction.

This first slide shows an example of the kinds of changes I'm talking about. These data come from a simple animal learning situation in which a rat is placed in a chamber with a hole on the wall, and occasionally food can appear in that hole, the so-called food magazine. Food is signaled by a diffuse auditory or visual event. As the animal learns that the signal means the coming of the food, the animal increasingly goes over to the hole on the wall in anticipation of the food and sticks his head in and out of that food magazine. If you put a little photocell in front of the food magazine, as we have, then you can detect his putting his head in and out during the signal before the food comes. That's what we use as a measure of his learning to anticipate food when we've arranged a relationship. As you can see at the left, when you arrange a relationship repeatedly, initially there is an increase in number of entries into the food magazine, and then when you destroy the relationship, as you do in extinction, there is a decrease.

Now, it's common to attribute these changes in the animal's behavior to changes in underlying associative learning on the part of the animal. We've learned over the last 50 years an amazing amount about these kinds of learning processes. We have wonderful databases and we have elaborate theories, which do a pretty good job of dealing with most of the data at a behavioral level. But there remain a great many, very basic questions about the nature of associations that we've been unable to answer -- questions we can't answer largely because the answers demand that we are able to compare the size of an associative change in one stimulus with that in another when the two stimuli are at very different stages of learning. We can't make that comparison because our theories of associative learning have been very elaborate, but our theories of the mapping of that associative learning into the behavioral indexes have been highly primitive. In fact, we've mainly assumed that whatever that mapping is, it simply preserves an ordering. It only preserves

monotonicity such that if one stimulus has more association strength than another, it will elicit more behavior. We've been reluctant to make stronger assumptions than that. But unless you do make stronger assumptions than that, there're a lot of questions that are basic questions you can't answer.

Let me show you an example of what I mean. Perhaps the most obvious feature of this curve is what you might call the negative acceleration of the change. That is, the changes that occur in acquisition are rapid at first and then decrease. Similarly, the changes that occur in extinction are rapid at first and then decrease. It's been common to attribute that negative acceleration in the behavior to a negative acceleration in the underlying learning process. That is to say, the association changes in the same way. But unless you're willing to make stronger assumptions about the mapping of learning into performance, the fact that the behavior is negatively accelerated does not imply necessarily that the learning is negatively accelerated. You could perfectly well have learning which goes on in a linear fashion, but the behavior may be limited physically by how often the animal can stick his head in and out of the magazine. Or it's limited on the extinction side because the behavior has got to go to 0. Wherever behavior starts, it's going to come down to 0 and it's going to look negatively accelerated even if the underlying learning process is linear. But unless you make stronger assumptions about the mapping of learning into performance, it's hard to see how you're going to infer that associative change is negatively accelerated.

Of course that hasn't prevented theorists for making the inference; that is, it's regular to assume in theories that learning is negatively accelerated. That was true of the early learning theories of, say, Hull and Tolman and Guthrie. It was true of so-called error correction models, which I'll talk more about today, models in which the organism is seen as comparing his current state of learning with some asymptotic level and on every trial changing by the amount of that error or that difference. That will naturally produce a negatively accelerated curve. It's true of almost all connectionist models. Theorists have been happy to make this assumption. But it would be nice to have something more than the word of theorists, because we all know that's not worth a great deal. We thought it would be useful to see if there were some other techniques we can devise that would help us to address this question.

Let me try to show you a way in which you can try to answer this question without strengthening your mapping assumptions. The next figure shows the basic idea of a negatively accelerated curve: if I give a few trials of training to stimulus which is down low on the curve (B), I'll get a big increase. However, if I give the same few trials to a stimulus (A) which is high up on the curve, I'll get a smaller increase. That's the basic idea of negative acceleration. If you do that, of course you can't compare the associative increase

you get in B and A unless you have stronger assumptions about mapping or unless you devise some other technique. Here is the other technique, which we've tried.

Let's suppose we have four stimuli, not only A and B, but A, B, C and D. Let's suppose in phase 1, we train A and C just a moderate amount. They're not untrained, they're not fully trained, somewhere in the middle. But we leave B and D untrained. Now the question is this. If I give the same number of trials to A in phase 2 as I give to B in phase 2, which stimulus will undergo greater change? If the learning curve is negatively accelerated, A will undergo less change because it's already high on the curve whereas B will undergo more change. But of course I can't just look at them because they start out at different levels and I don't know how to compare them. But I can look not at the individual stimuli, but at the compounds.

Now, let's think about why looking at the compounds is going to help. Let's suppose we hadn't done phase 2; we just did phase 1 and tested. If we did that, AD and BC should be equal because each of those compounds contains one stimulus that was trained, A or C, and one stimulus was not, B or D. In fact, if we've counterbalanced our stimuli across animals, those two compounds are the same as each other. They have to equal each other. If we hadn't done phase 2, those two compounds have to be the same. But we did do phase 2. When we do phase 2, that can produce some differences between responding to the compounds, but the differences all have to do to what we did in phase 2. For instance, if we get more associative change in B than in A, then we should see more responding to BC than to AD, because the increase in BC would be bigger than the increase in AD due to the fact that the increase in B was greater than the increase in A. But BC and AD in the absence of that training started exactly at the same level, so we can now compare the magnitude of the changes because they're occurring from the same level. That's what we did in this very first experiment. We did magazine training with four stimuli, two actually trained in phase 1 and two not. Then in phase 2, we did some more training of A and B, the same number of trials for A and B, and then we tested.

The next figure shows what happened in that second phase, when we're now training B for the first time and A gets continued training. I have shown responding in the pre-stimulus period before any stimulus comes on, just to give you an idea of how much he is sticking his head in there anyway. You can see that A starts out high, it's been pre-trained, and it goes up some, not surprisingly. B starts out, of course, low and it appears to go up more. It looks like we've got more change in B than A. But of course this is the problem, we don't know how to make that comparison, because A might be stuck on a ceiling. We're going to now test the compound.

When we test the compound, you see in the next figure that we get more responding to BC than to AD. That must have mean we've got more learning in B from a fixed number of

trials than we did in A from that same fixed number of trials. That is to say, we've got negatively accelerated changes because we've got more learning in B than A. We can make that comparison by looking at the compounds in a way which doesn't require any strong assumptions about the mapping of learning into performance.

Now, you're going to say "well now why do you go through this long song and dance? I thought I knew learning was negatively accelerated. You told me I didn't know learning was negatively accelerated and then you told me it's negatively accelerated after all. What's the point of it all?" Well, I think it's very important in science to make a distinction between those things you actually know and those things which you only think you know. If you look in the history of discussion of science, you'll see I'm not of course the first who thought of this. So for instance, just to pick an eastern scholar, Confucius said "real knowledge is to know the extent of one's ignorance." Similarly, Hazlitt, who is one of my favorite scholars from the 18th century in England, said "the origin of false science is the unwillingness to acknowledge our ignorance." I think that's right. A lot of work I have done over the years has in fact had the structure in which we think we know something, but when we look at it carefully we don't. So we ask the question, what does it really take to know it and try to answer that question. In this case, it turns out we were right, negative acceleration is occurring in the underlying associative process in acquisition. But we needed a better technique to in fact find that answer.

But of course the same question arises in extinction. It's useful to ask how would we address the question in extinction. The next figure shows the issue in extinction. Extinction is negatively accelerated if the changes you get from, say, two trials early in extinction are greater than the changes you get from those same two trials later in extinction. Let's ask how could we address whether associations show negatively accelerated change in extinction that using this compound test.

Well, here is the one way. It's a little more complicated experiment because you have to train before you can extinguish. It starts out training four stimuli, A, B, C, and D; two are auditory, two are visual. They're trained up essentially to asymptote and then two stimuli, A and C, get a moderate amount of extinction, enough to bring them down some, but not enough to take them all the way down. Now, the question is, if I give the same number of further extinction trials to A, which has been pre-extinguished and B which was not, which will undergo greater loss? The logic is really the same. If I don't do that extinction in phase 2, then these two compounds have to be the same. They each contain one stimulus which has been moderately extinguished and one of which has not. And again with counterbalancing, they really are the same across different animals. So, there is no way they can be different. If I see a difference after extinction phase 2, the difference has to be due to

what I did in the extinction phase 2. That's going to tell me which one of those underwent greater loss. That's exactly what we did, again with that magazine approach situation.

The figure shows what happens in the test. Pre-period, things are very low. AD shows a moderate amount of responding whereas BC shows a lot less. That means that the decrease we've got from those extinction trials on B must have been greater than the decrease we've got on A, which is to say the decrease early in extinction was greater than the decrease later in extinction in the association. Again, just about any performance assumption we make which preserves monotonic mapping is going to be able to use this kind of test to make the inference that the extinction was also negatively accelerated. This again is something we've looked at in lots of preparations; I've just shown you the example of the magazine approach because it's convenient.

Now, early on, I said the one reason people thought learning might be negatively accelerated is due to error correction models, models which say the organism compares its current state with some end state, treats that as an error and changes some constant amount of that. A primitive early example of an error correction model is the Rescorla-Wagner model, which I've shown you in the next figure. Let me just walk you through that model, because we're going to need some of its implications.

The figure shows what happens according to that model on two kinds of trials, trials upon which stimulus A is followed by US and trials upon which a stimulus compound of A and B together is followed by US. The expressions simply tell you the amount of change in association that is supposed to happen when you do that. The change is in associative strength, V_A - this theory uses V_A as the sign for association. Change in V_A is a constant, K , proportion of the difference between the current associative strength, V_A , and some asymptote, λ . The organism views its current associative strength, looks at the difference from its asymptote, and changes some constant amount of that difference. As it does that, of course, the difference gets smaller over trials. So, you'll get a negatively accelerated acquisition curve. I mean that's how the theory is constructed.

The simple equation for A-US trials is a very old idea really. What was new in this model were assumptions that were made about the AB trials. The same basic idea applies there, the animal is calculating an error and correcting himself. But what's interesting in this version of this kind of model is that when the animal identifies the error, he takes into account the current strength not only of the stimulus which is changing but also the other stimulus which is there on the trial. The change in A depends upon the current strength of both A and B. The current strength of the AB compound is calculated and used as the basis of the error. In this model, the strength is the linear sum of the elements, this is not really an essential feature.

This common error is what allows this model to explain modern phenomena, which have been so important in conditioning, like blocking. Let me remind you what blocking is because blocking is probably the most important thing to happen in conditioning in the last 50 years. If you pre-train one stimulus called A to asymptote and then present AB followed by the same US, that B stimulus, despite the fact it is followed by US, will be blocked in its conditioning by the prior training of A. Prior training of A enables A to prevent conditioning of B if A accompanies B when it's reinforced. If A hadn't been there, B would get lots of conditioning, but with A there, conditioning is blocked. This model says the reason that happens is that the animal is trying to calculate the error, but is using both A and B. When I pre-trained A, A is already high; then the sum of A and B is already high, so there'll be no change in B on those trials because there is no error. Models like this have been highly successful doing a great job not only predicting old things, but generating new predictions. But such models, despite the fact they're very popular and extend to connectionist networks as well, make a very strong prediction which has not been tested. The strong prediction is that on a trial where A and B are followed by the same US, A and B will change in the same direction by the same amount. No matter how different A and B are before the trial, one can be very high and one can be very low. But if they're both reinforced together, they both have to change the same way and by the same amount. It's a very strong prediction; moreover, it seems very unlikely to be true. All these models make that prediction, but it's been impossible to test. In order to test it, you have to compare the associative change you get in A, which is already high with the associated change you get in B, which is very low.

But using the compound technique I've been talking about, you can ask that question and you can decide whether the strong prediction is right or not. I want to show you an example of several experiments which have done that because they turn out to have interesting different features. The first one was a design which was intended to provide a very strong test. We decided we take an A and a B, we reinforce them together, but we would give them as different beginning associative strengths as we could manage. We used what's called a conditioned inhibition paradigm. That's pre-training where A gets reinforced on its own, making it very strong, but not reinforced when B is there, making B a very strong inhibitor, because B predicts that A will not be reinforced. You end up with an A, which is very excitatory: and a B, which is very inhibitory. We did that for two pairs AB and CD.

Then we ask what would happen if we take the AB compound, which has a very high element and a very low element and we reinforced the compound. Will these two elements change the same amount as these theories all say they should? Now if you think about it, I think you can convince yourself any of three outcomes is reasonable. It could be that when I reinforced AB, A will change more than B, because after all the organism has already

identified A as a good predictor. An attentional theorist might expect the animal to blame A for any increase, so A would increase more than B. Alternatively, you might expect B would undergo bigger change, because after all B as an individual stimulus is very far from the asymptote. If I'd reinforced them separately, A on some trials and B on other trials, surely B would change more than A. Or you might think they'll change by exactly the same amount, which is what error correction models say. The claim I'm making is we separate these alternatives by testing compounds, so let's think about the test.

Let's suppose we hadn't done this compound conditioning phase. When we now present AD and BC, they ought to be equal because each of them contains one stimulus which is excitatory A or C; and one stimulus, which is inhibitory B or D. Although the B and D were made inhibitory with a different exciter, they should equally transfer their inhibition to the other stimulus. So if I hadn't done this conditioning phase, those two compounds would be the same. Now, we can ask are they the same when we test them after doing that conditioning? If it turns out, for instance, that BC is bigger than AD, that's got to mean B picked up more than A on those trials. If it turns out that AD is bigger than BC, then A must have picked up more. If they're still equal, that means they picked up the same amount. The test compounds are starting from the very same level, so we don't have to worry about the levels problem.

Here is what happens when you test them. When you test them, you find out the animal shows responding to AD and responding to BC, but the responding is much greater for BC. What that means is that B must have gained more on those AB trials than did A. That means the inhibitory stimulus, B, showed more increase than the excitatory stimulus, A. They don't change by the same amount. All the theories which assume they're going to change by the same amount are wrong. We couldn't have known that without this kind of compound test. Now, this is very important because it not only bad for the Rescorla-Wagner model, which all right, we can live without that, but it's bad for a whole host of other models including a lot of error correction network models, which all make this assumption that they change by the same amount. If they don't change by the same amount, those models don't work. This is not a trivial assumption that those models made, but essential for them to work.

Now, this is one of many experiments we've done like this. We've used not only an exciter and inhibitor, but we've used other stimuli like an exciter and a neutral stimulus. We've used other conditioning preparations. And we've looked at extinction, in which the one that changes more is the excitatory stimulus, not the inhibitor.

Now one variation that is of particular interest is a variation in which you don't simply retain the same reinforcer from stage 1 to stage 2, but you change the size of the reinforcer.

Now, I want to tell you about an experiment which does that, because it provides us some additional information. Let's suppose we repeat the same experiment except that the reinforcer we use in phase 1 is stronger, which I've indicated in the figure by a double plus, than the reinforcer we use in stage 2. Now, think about stage 2. When we follow the AB compound with the reduced reinforcer in stage 2, then from the point of view of the B stimulus, that's an increase in reinforcement. But from the point of view of the A stimulus, that's a decrease in the reinforcement. If I had trained in phase 1 and then trained A with a single reinforcer and B with a single reinforcer, they would become closer together. They would both move towards the middle. The A stimulus would come down because it's experiencing a decrease in the reinforcement value and the B would come up. Any error correction model which views the error as based on the differences of individual stimuli from the asymptote anticipates that reinforcing the AB compound should produce convergence of A and B. But an error correction model that says you have a common error term, like the Rescorla-Wagner model, says that despite the fact that A is now being followed by a weaker reinforcer than before, it will go up. It will go up because the AB compound is under-predicting the reinforcer. There is still an error because the B is inhibitory keeping the prediction by AB down. The Rescorla-Wagner model says A and B will both go up. They'll both go up by the same amount. Individual error correction models say no, A will come down and B will come up. The same test as before can be used to decide.

In this case, it's also useful to test the individual stimuli. Notice that, for instance, we can compare A and C. One of those two stimuli, A, got this additional training and C did not. By looking at A and C, we can ask did A go up or down because A is being trained and C is a control stimulus. Similarly for B and D, we can ask did B go up or down? Of course, we can't figure out which one went up or down more, but that's why we've got the compound test.

Let's look at the individual elements first. The Rescorla-Wagner model says A and B will both go up, and it turns out Rescorla-Wagner model is right. B went up and A went up. Notice A went up even though it's now being followed by a weaker reinforcer than it had before, because it's being followed by that weaker reinforcer in compound with B. One point for the Rescorla-Wagner model and ones like it. But we can't tell which one went up more. It actually kind of looks like maybe A went up more, but it's hard to make tell because of their different parts on the scale. To tell which one went up more, we need to look at the compounds. The result is not so good for the Rescorla-Wagner model. Comparing responding to BC and AD looks like the result of the previous experiment. Because BC is more responded to than AD, it looks like B went up more than A. In summary, when you do this experiment, both stimuli go up. They go up together; they don't converge. But they go up by

different amounts. It's not good for anybody. It's not good for single term error correction models because both stimuli go up. It's not good for the compound error correction models because they don't change by the same amount. So there is a real problem here for any class of error correction models.

Again, this is just one example of a lot of experiments we've done like it. We've done it in autoshaping and flavor aversion and various other kinds of situations. The compound technique, which allows us to make this assessment, has really led us to make serious inroads into addressing an important theoretical question that we couldn't have answered before.

Moreover, it turns out the compound technique will allow you to answer lots of other general questions we need to answer about associative learning. I want to tell you about some other cases where we've applied it to other important situations.

One question which I've always thought was interesting is whether it matters what the initial value of a signaling stimulus is before you do conditioning? Now, in fact if you look at conditioning experiments, people go to all sorts of lengths not to use CSs which have initial value. In fact, if you read traditional descriptions of Pavlovian conditioning, they'll say in conditioning you take a neutral stimulus and follow it by some important stimulus and the neutral stimulus changes. We have almost no experiments in which we don't take a neutral stimulus. We don't ask whether the difference in the initial value of the stimulus matters; but that's an important question about the nature of association. Is the learning the same regardless of the initial value or does the initial value make a difference?

Now, of course, in traditional conditioning preparations, most of our CSs have been carefully picked to be neutral, so we have to find one in which they're not neutral if we want to study this. It turns out that the flavor aversion preparation is one in which you have natural stimuli for doing this. Some years ago, we did a bunch of experiments with four flavors; sugar (sweet), salt (salty), quinine (bitter) and mild hydrochloric acid (sour). Now, it turns out that with some mixing of the different concentrations, you can arrange it so that two of those stimuli are mildly positive: sweet and salt. You can match them against water and get them to the same level of positivity in a choice situation. Two of them are mildly negative: bitter; and sour. You can again match them for their level of negativity. Now, you can ask what would happen if we try to condition a negative stimulus and a positive stimulus using a negative consequence like lithium chloride. Lithium chloride is a commonly used reinforcer that makes the animal feel sick and makes him reject paired substances. You can use those stimuli and that kind of conditioned US to answer this question of what difference initial value makes.

Now again you can imagine three plausible answers. If I have a stimulus which is

inherently negative and I try to condition it to be negative, will it be especially good at that? A lot of clinical psychologists seem to think it would. If you think back to the example of phobias where spiders are always pointed out as common phobic objects and flowers are not, the assumption is they've got some initial negativity and that makes them specially susceptible to negative conditioning. That is, you might think that an initially negative stimulus would be easy to condition negatively. Or you might think that it will be hard to condition because that's already part way there. If it had been made negative by some conditioning operation, it would be like it's on the negatively accelerated acquisition curve. It could be the unconditioned negativity would function just like conditioned negativity and make it undergo less change. Or the initial value might not matter. It could be that the positive and negative initial values would not matter; the amount of change you get might be the same, and it just builds on top of the initial difference.

Here is a way to try to address it. For this experiment, we have two positive and two negative stimuli. I've called them P1 and P2 and so forth. The P1 and P2 are sugar and salt counterbalanced and the N1 and N2 are hydrochloric acid and quinine counterbalanced. The same animal now gets conditioned to one positive P1 and one negative N1. Now, we ask what does it do to the compounds? Again let's go back and imagine we hadn't done the conditioning. If we hadn't done the conditioning and we presented those two compounds, P1N2 and P2N1, they ought to be equal. Each of them contains one inherently positive and one inherently negative stimulus. If they're counterbalanced, there is no way they could be different. If we get a difference in the test between those two compounds, if the first compound P1N2 is greater than the second P2N1, that means that the positive stimulus, P1, must have undergone greater change than the negative stimulus, N1.

The first time we did the experiment, we did it with lithium chloride, conditioning of positive and negative stimuli with a negative consequence. What I plotted here is amount of conditioning. What you can see here is we got more conditioning when we trained the positive stimulus than when we trained the negative stimulus. This is with aversive training.

Now, unlike the other situations we've been talking about, these positive and negative stimuli cannot necessarily be thought to be equally salient. They cannot be counterbalanced because they're inherently positive and negative. With our other stimuli, they could be counterbalanced, so we didn't have to worry about this. But here it could just be that our positive stimuli are easy to train or negative stimuli are not easy to train. That is, the result might have nothing to do with the fact that we trained them with aversive stimuli. Consequently, it's important to do the experiment not only with an aversive consequent, but with a positive consequent. It turns out you can do that too. You can pair these stimuli not with lithium chloride but with a strong concentration of sucrose and Polycose. That

concentration is something which rats love. I can make something into a positive stimulus by pairing it with this positive compound. Now, I can do the same experiment, but use a positive consequent, again conditioning P1 and N1. The figure shows that with the appetitive consequent, I get the opposite result. With the appetitive consequent now the negative stimulus shows more conditioning compared to the positive stimulus. What this means is that if I have an inherently attractive stimulus, it will be easier to make negative or easier to increase its negativity than if I've inherently negative stimulus. The inherent value is acting like a conditioned value. The stimulus will be easy to change or will change a lot if there is a very large discrepancy between its inherent value and the consequent. The spider-flower story that's told by clinicians is just backwards here.

Now, related to inherent value are several issues about acquired value and we might ask questions about how easy is it to train various stimuli which have various kinds of acquired states. One question, which is related to this one, is how hard is it to retrain a stimulus after extinction? If you read the conditioning literature or for that matter any introductory textbook, they'll tell you that a stimulus which has been trained and extinguished is easy to retrain – easier to train than a stimulus that's not been trained and extinguished. Retraining is faster than initial training. But if you think about that, you realize there is an inherent ambiguity in that claim. Is retraining faster because the stimulus which is undergoing retraining undergoes fast associative change or is it faster because the stimulus which is undergoing retraining was never really fully extinguished and simply has a head start. Is the retraining stimulus smart so it learns quickly or has it a head start because you never took it back to 0? They're very different. We can ask the question is it really true at the associative level retraining is faster than initial training or is that an accident of our behavioral measures?

Here is a way to ask that question. Take two stimuli, A and C, train them up to asymptote, fully trained. Now extinguish A and C, again fully extinguishing them. Two other stimuli, B and D are available, but receive no training. Now, retrain A, that's retraining, and train B, this stimulus for the first time. That's the usual comparison: how fast will A come up, how fast will B come up; typically A comes up faster. But you don't really know if that means there was greater associative change in A or in B; all you know is A had more behavior. So, in order to make the comparison, we have to do something like a compound test. Again imagine we didn't do the retraining. When we test AD, we have the compound, which has one stimulus that's just been extinguished, A, and one which has not been trained, D. When we test BC, also have a compound which has one stimulus which has been extinguished, C, and one that has not been trained, B. If we haven't done retraining, they're the same. Now, the question is did retraining produce more associative change in the

previously extinguished stimulus or the previously untrained stimulus? We can ask do we get more associative change when we retrained or when we did initial acquisition?

The one figure shows what happens in the usual test when we simply look at how A and B change and you can see that A comes up very rapidly and B comes up much less rapidly. Those are the data that are usually used to say retraining is faster than initial training. But you'll notice there is this awkward fact that they're little different to begin with. You can be sure you're always going to see that if you look carefully enough. So it's really a lousy comparison. Rather, we have to ask about the compound. It turns out when you test the compound, you get more responding to BC. That means you've got more increase to B than you did to A. When you did retraining, although behaviorally the A stimulus went up faster, in fact it was learning less. The neutral stimulus is learning more. The reason I think is pretty clear. It's because the faster retraining is due to the fact that the extinction was just never fully complete. It looks like faster retraining and it is in terms of behavior, but it's not faster learning. It's faster showing of learning that was already there before. What this means is that fast retraining is due to incomplete extinction, not due to the fact that the stimulus actually acquires the association more rapidly after extinction. This is a very important difference that you need to know the answer to.

Now, there is a very similar case where you don't look at retraining, comparing a retrained stimulus with a neutral stimulus, but you compare the training of the conditioned inhibitor with the neutral stimulus. A standard way of assessing whether stimulus is a conditioned inhibitor or not is to try to train it as an exciter and compare how fast it trains compared with how fast a neutral stimulus trains. The faster behavior to a neutral stimulus is taken as evidence that the inhibitor is inhibitory. But here is the question. Did the inhibitor acquire its behavior more slowly because it started at such a severe disadvantage or did it acquire more slowly because it's hard to turn the inhibitor around? It's hard to convince the stimulus now that it ought to have the opposite value; that is, the associative learning is slow. I think again we tend to confuse these in our discussions, but they're very importantly different. Is an inhibitor slow to become an exciter or is it becoming an exciter fast but just starts from so low we don't see it?

Again, you have to have the compound test to answer that. Here is our experiment. In this experiment, we used five stimuli. One stimulus I've called an X. It's simply an exciter that's used to train two inhibitors. We train X as an exciter and A and C as inhibitors and we leave B and D as neutral. Now, we ask who trains faster: A, an inhibitor; or B, a neutral stimulus. We're going to compare the compounds too.

Let's look first at phase 2, which is where the usual data come from. In phase 2, I've shown you how responding develops in the pre-stimulus period, in the A stimulus, which

starts out as an inhibitor and the B stimulus, which starts out neutral. You can see behavior develops more rapidly in B than in A. That's the standard evidence from the so-called retardation test of A as an inhibitor -- A comes up more slowly than a neutral stimulus. But of course we don't know did A come up more slowly because it was having trouble learning or did A come up more slowly because it started at such a disadvantage?

Looking at the compound tells you the answer. In the compound, we can ask who showed more associative increase, A or B? It turns out that AD shows more responding than BC. That means A, the inhibitor, underwent greater associative change than B, the neutral stimulus. A in fact learned faster. It just learned from such a disadvantage that it had a long way to go. That's not really surprising if you believe in error correction models. In error correction models, which has the bigger error, a neutral stimulus which is far from the asymptote or an inhibitor which is even farther? The inhibitor on each trial is in fact learning more, but it just has trouble catching up. If all you want to just decide is whether or not you have an inhibitor it doesn't matter. But if you want to understand what's going on in the nature of the process, if you want to measure which stimulus is undergoing more associative change, you have to use something like the compound test. It gives a kind of surprising answer that the inhibitor is in fact undergoing faster associative learning.

Another example: this is an example which is close to the heart of theorists. One question which has been of interest to theorists: is the learning process faster in acquisition or in extinction? Do you undergo acquisition faster or slower than you undergo extinction? Now, if you actually look at behavior in the kind of graph I showed you before, it sure looks like the changes are faster in extinction. This is pretty typical. But if you ask theorists to tell you about what is sometimes called the rate parameter, how fast does learning happen and how fast does extinction happen, they'll tell you acquisition happens faster than extinction, despite those data. The question is are they right? Are theorists right that in fact acquisition is faster than extinction in the face of data like those?

Now, that's just not a casual decision that theorists are making. Many theories, in order to explain certain phenomena, require that acquisition be faster than extinction. If it's not, certain phenomena simply can't be predicted. So it's an important question, but it's really hard to see how you're going to answer that question. I mean you're asking: if I give three trials to a stimulus that is low and which I'm asking to go up, and three trials to stimulus that is high and which I'm asking to go down, which changes more? You've got not only the problem that they're changing from different levels, but also that they're changing in different directions. How can you ever ask which one really was faster? Well, it turns out -- it won't surprise you -- that the compound test allows you to answer that. It takes a little different variation on the test, but it does address it.

The figure shows how it goes. Four stimuli: A and C are trained while B and D are not. Now think of phase 2 as reversing the A/B discrimination, but not the C/D discrimination. What we've done, we train A and C, not B and D and then we reverse the training on A and B, but keep constant the training on C and D.

Now, of course if we hadn't done this reversal phase and we now tested AB, it would be equal to CD, because AB contains one trained and one non-trained stimulus and CD contains one of each. But when we do carry out the phase 2 reversal, if we test AB versus CD, we can ask what changed more, A or B? If A went down more than B came up, then AB would be less than CD. If A went down less than B came up, then AB would be greater than CD. We can ask who changed more A or B, because if we hadn't changed them, they would be equal to CD.

The figure shows what happens during the second phase. It's kind of cluttered graph, but the dots show where things were before we engaged in phase 2. At the end of phase 1, two stimuli are high, two are low, and there is a pre-stimulus. Now, in the second phase, C is just treated the same as it always was and stayed up. D is treated the same as it always was and stayed down. But A is extinguished and B is conditioned. They both change in the right direction. The question is, is the underlying associative change, the decrease in A greater or less than the increase in B? Now, it kind of looks like the decrease in A was greater than the increase in B, but of course the problem is we can't compare them. The only way to compare them is to look at the compounds. If in fact A has gone down faster than B has come up, then AB will be less than CD. But that's not what happened. In fact, AB is greater than CD. AB being greater than CD means that A underwent less change than B. The B went up in acquisition more than A went down in extinction when both received the same number of trials. The theorists are right that in fact acquisition is faster than extinction. We can actually make the comparison for the first time.

Now, I have one final experiment to confuse you with. In the last experiment, the reason that we trained C and D in the same old way during phase 2 was that we were worried that what we did to A and B would generalize to C and D. We wanted to hold C and D where they were by continuing to treat them the same way, so we didn't have to worry about any generalization from the other treatment. But it turns out that if you don't hold C and D the same, that changes the experiment slightly, and you can ask a related question. The related question is this: When I train a stimulus, it will generalize some amount to other stimuli that are similar and when I extinguish a stimulus, it will generalize too; but which generalizes more, acquisition or extinction? It turns out that's been an historically very important question, but it's really hard to know how we're going to answer that question because we can't compare changes at different levels in different directions. How do I

compare the widths of these two generalization gradients? Well, people have been happy to do it, but they shouldn't have.

It turns out that a small variation on this experiment will allow us to answer that question. The figure shows the design. It looks very much like the last one except I have changed the notation to emphasize the stimulus similarities. I've got two stimulus dimensions, A and B. In this example, the two exemplars of A, A1 and A2, are both trained. The two exemplars of the B dimension, B1 and B2, are both not trained. Then we extinguish A1 and we train B1 and then we test these compounds.

Let me walk you through the logic of why we're doing this. The logic turns out not to be that easy to see unless you do it step by step. I can never remember it. Let me show you how the logic goes. Let's suppose that we find that A1B1 is greater than A2B2. Let's suppose that's the result when we test these. That suggests in turn that the difference between B1 and B2 is greater than the difference between A2 and A1. You can see this by a simple algebraic manipulation: transpose B2 over to left side of the equation and A1 to the right side, changing signs appropriately. If the first equation is true, then the second is going to be true. But if the second equation is true, that says that difference between B stimuli, the treated stimulus, B1; and the generalized stimulus, B2, is greater than the difference in the As between the treated and the generalized stimulus. The Bs are further apart than the As, if that's true. That means the B1 has generalized less than A1. Because B1 is more different from B2, it must have generalized less than A1 and A2. That in turn means that the reinforcement generalized less than non-reinforcement. Now, you may have to write it down and go over in your head five or six times. But the logic is right that if we find that the A1B1 compound is greater than the A2B2 compound, that's telling us that reinforcement generalizes less than non-reinforcement. For the first time, we can really answer that question in a meaningful way.

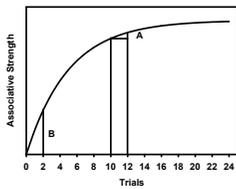
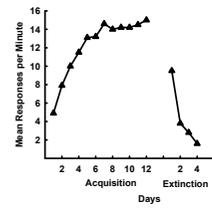
Now, the next figure shows the data from phase 2 demonstrating that in fact in training, the extinction and the acquisition were complete. This matters because these changes happen at different rates and you don't want those rates confounding things; you want both of these changes to be completely done before the critical training phase. The next figure shows the test. The test result was A1B1 is greater than A1B2. I'm sure you all remember what that means. That means that reinforcement generalized less than non-reinforcement. In fact, it's true that the breadth of generalization is greater when you're extinguished than when you've reinforced. It turns out theoretically to be a very important issue.

These then are a whole set of examples of how you can use this compound technique to answer a bunch of questions. Let me review with you exactly how we have tried to use it to answer various questions. Some were questions we thought we knew the answer to but

didn't, some were questions we thought we could not answer, some were questions we did not think to ask, and some were questions on which we are just beginning to get the answers. The first question we looked at was the shape of the learning curve. Now, here we thought we knew the answer. We thought it was negatively accelerated. It turns out that if you tested in the right way, the answer you thought you knew was in fact right. But you had to test in the right way. You didn't really know the answer before. Then we've used it to try to evaluate important error correction models. For the first time I think we've provided decisive problems for major error correction models. Not that the Rescorla-Wagner model was without problems before, but this is a really major problem because this is an inherent assumption it can't avoid and it turns out to be wrong. Things do not change by the same amount. But we couldn't know that until we had the compound test. Then we ask what would happen with conditioning with valued stimuli? It is the question which basically we've been reluctant to ask before, but now we can ask it and we can get an answer. The answer is that you get an amount of conditioning which was the same as if the inherently valued were conditioned instead of inherently valued. We can also ask about stimuli at different starting points such as conditioned inhibitor and neutral stimulus and such as a retrained stimulus. There we were able to make an important distinction between two stimuli being different because they started at different places or because they changed at different rates. We asked about changes in a different direction where we can ask is the decrease in extinction the same or different from the increase in acquisition. Finally, we asked about stimulus generalization where I think the technique has shown some promise.

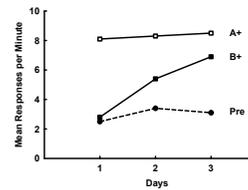
This technique is not without its problems. It's not a perfect solution. Sad to say there never is. For instance, it assumes that compounds don't introduce anything which is not in the elements. That is, there is not some new process which occurs when you put things together in compound. That's a very common assumption which most people make, but it is an assumption. We've also assumed that we can somehow take the elements of a compound and add them together. The easiest way to think of is it's linear addition, but of course everybody does assume you can add them together somehow. It turns out it doesn't make much difference for this kind of a compound test by what rule you add them together. So, it seems to me this is a technique which has a lot of promise. It's already generated a lot of data which I think are of interest; and it shows a great deal of promise for generating new data. I hope some of you who are out there who are younger and not quitting the game like I am, will actually use the technique to some advantage. Thank you very much.

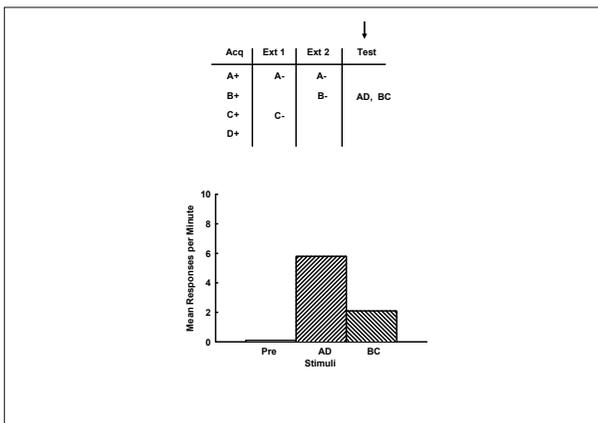
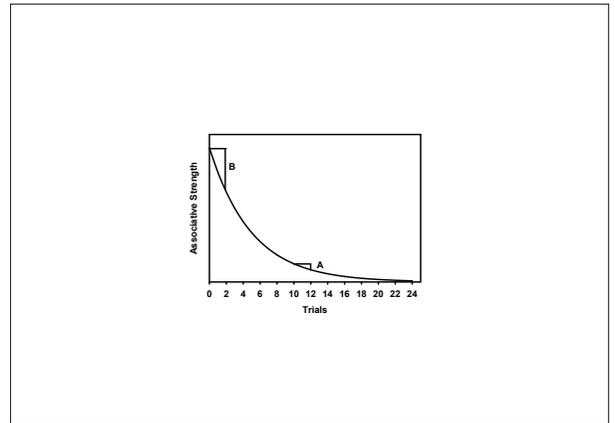
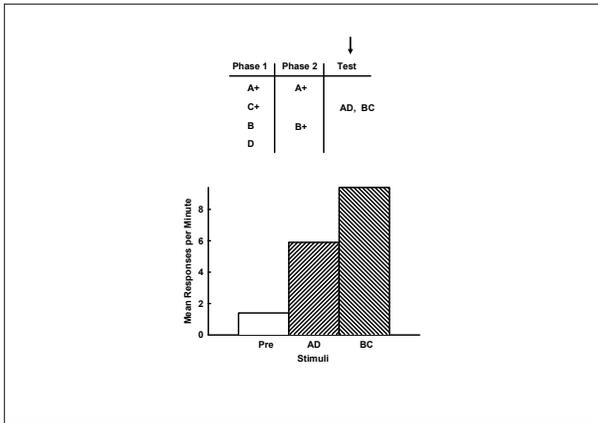
Measuring Changes in Associative Learning



↓

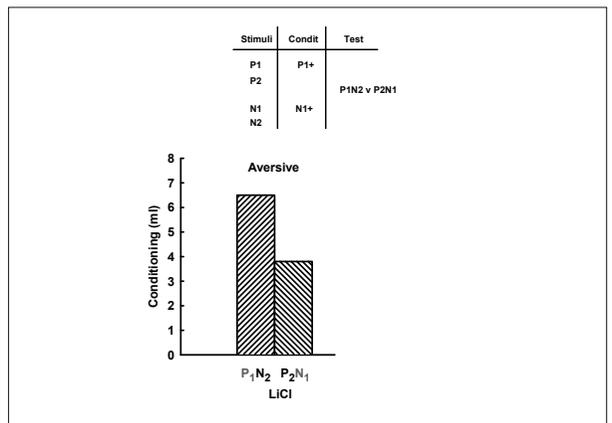
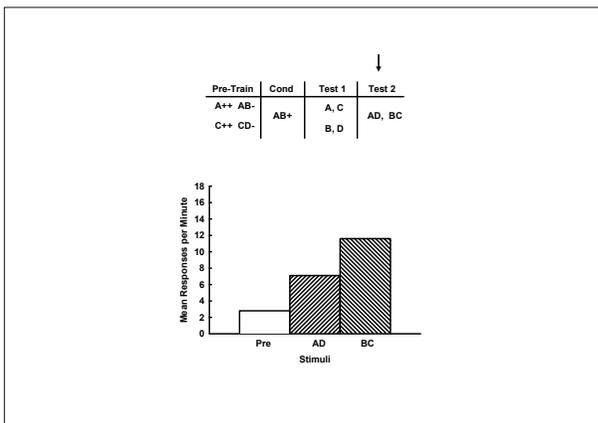
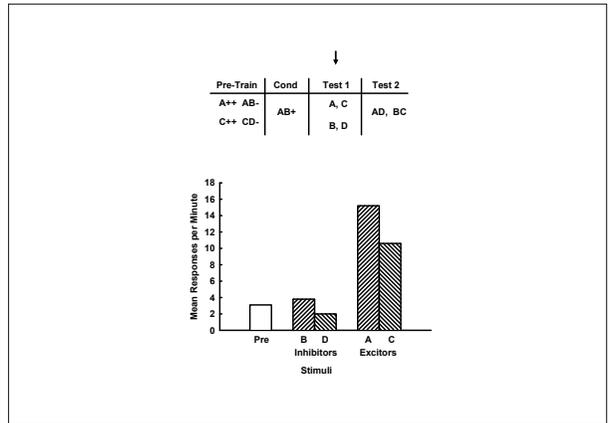
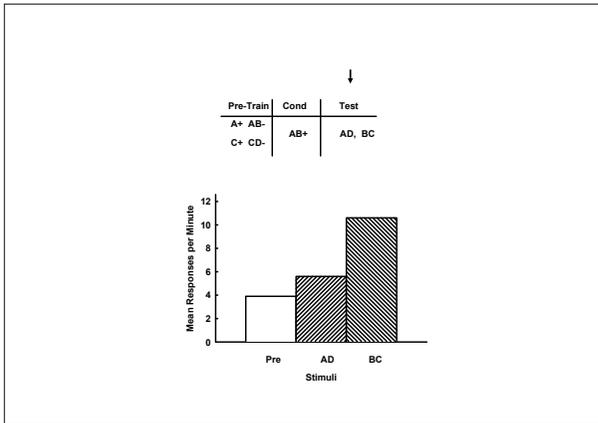
Phase 1	Phase 2	Test
A+	A+	AD, BC
C+		
B	B+	
D		

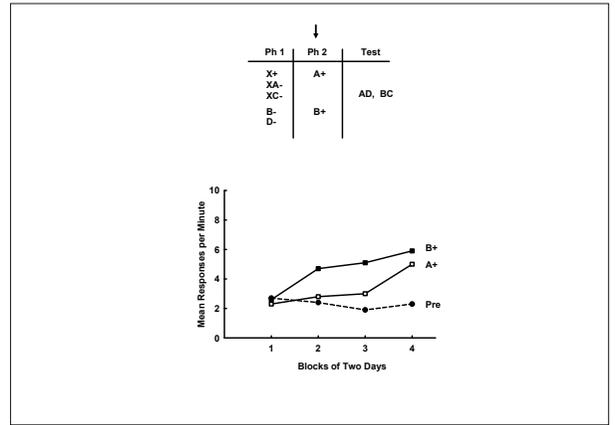
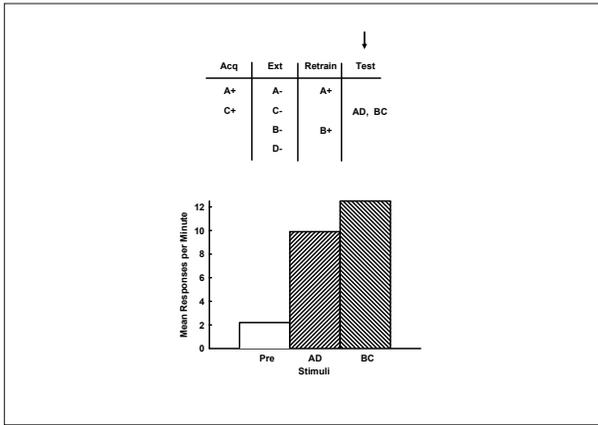
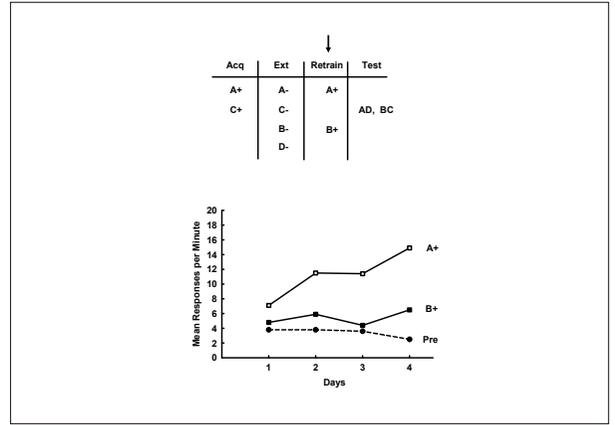
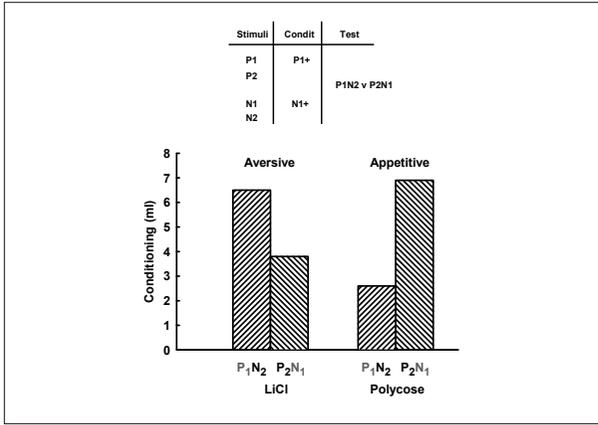


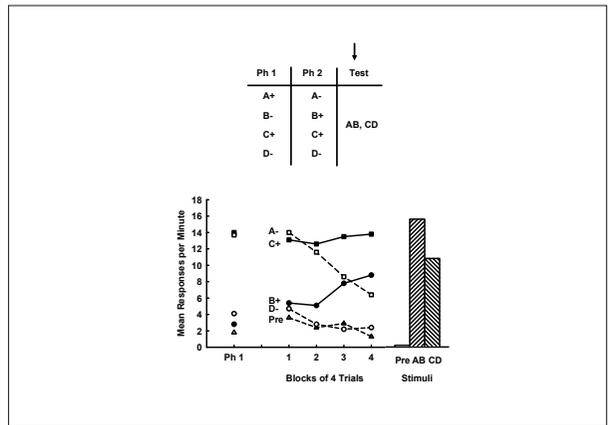
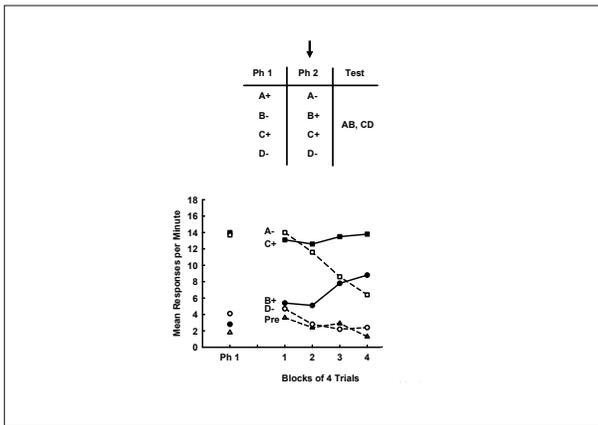
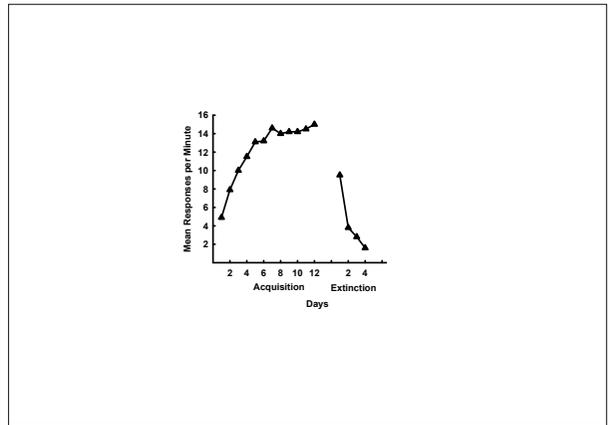
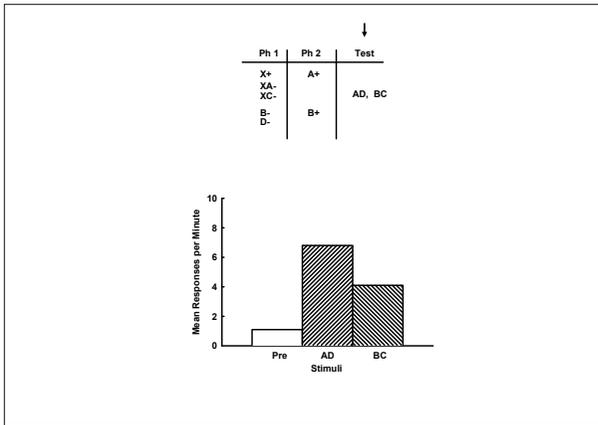


A -- US
 $\Delta V_A = k (\lambda - V_A)$

AB -- US
 $\Delta V_A = k (\lambda - V_{AB})$
 $\Delta V_B = k (\lambda - V_{AB})$
 $V_{AB} = V_A + V_B$





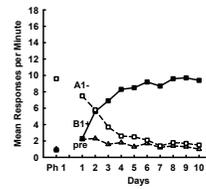


Ph1	Ph 2	Test
A1+	A1-	
B1-	B1+	A1B1
A2+		A2B2
B2-		

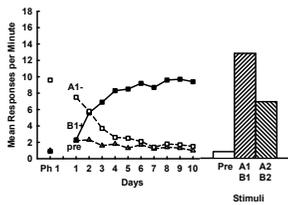
if $A1+B1 > A2+B2$ then
 $B1-B2 > A2-A1$
 $B1/B2$ difference $>$ $A2/A1$ difference
 B1 generalizes less than A1
 Reinforcement generalizes less than Nonreinforcement

↓

Ph1	Ph 2	Test
A1+	A1-	
B1-	B1+	A1B1
A2+		A2B2
B2-		



Ph1	Ph 2	Test
A1+	A1-	
B1-	B1+	A1B1
A2+		A2B2
B2-		



Compound Test Uses

- Learning Curve Shape
- Error Correction Models
- Conditioning Valued Stimuli
- Different Starting Points
- Different Directions of Change
- Stimulus Generalization