

# Elemental Causal Learning from Transitions

Kevin W. Soo (kevin.soo@pitt.edu)

Benjamin M. Rottman (rottman@pitt.edu)

Department of Psychology, University of Pittsburgh  
3939 O'Hara Street, Pittsburgh, PA 15260 USA

## Abstract

Much research on elemental causal learning has focused on how causal strength is learned from the states of variables. In longitudinal contexts, the way a cause and effect change over time can be informative of the underlying causal relationship. We propose a framework for inferring the causal strength from different observed transitions, and compare the predictions to existing models of causal induction. Subjects observe a cause and effect over time, updating their judgments of causal strength after observing different transitions. The results show that some transitions have an effect on causal strength judgments over and above states.

**Keywords:** causal learning; causal reasoning; time

## Introduction

Elemental causal learning is the process of learning whether a single potential cause has an influence on an effect. It is an ubiquitous mental process that helps us navigate and manipulate the world (e.g., Does premium gas get my car better mileage? Does doing a colleague a favor make them more friendly? Does watering my plant twice a week make it healthier than watering it once a week?) Consider the following extended example: a patient suffering from chronic fatigue tries a new drug that claims to boost energy levels. Over the next week, she takes the drug on three days (drug = 1, no drug = 0), keeping track of whether she has high (1) or low (0) energy. Figure 1 shows two possible patterns of experience with the drug, which we contrast.

	Day	Tue	Wed	Thu	Fri	Sat	Sun	Mon
(a) Drug		0	0	1	0	1	0	1
Energy		1	0	1	0	1	0	0
(b) Drug		0	0	0	0	1	1	1
Energy		1	0	0	0	1	1	0

Figure 1: Example Longitudinal Data Sets

From the experience in Figure 1a, how might the patient infer the drug's causal strength? Consider Wed-Sun. Between each of these days, stopping and starting the drug is accompanied by corresponding changes in energy, suggesting a strong positive causal relationship between the drug and energy. There are other days (Tuesday and Monday) that do not fit with this pattern, so the drug does not always work and it is not always needed. However, given the pattern from Wed-Sun, it is hard to explain the consistent pattern without inferring that the drug worked. It seems less likely that the drug has no influence on energy,

but by some coincidence some unknown factors changed the patient's energy at the same times (and in the same direction) that the patient happened to take the drug.

Contrast Figure 1a with Figure 1b. In Figure 1b, the three days that the patient took the medicine are all grouped together. Now it is much less convincing that the medicine works. Because they are grouped together it is more likely that the pattern is due to a coincidence; perhaps the patient had more energy on Saturday and Sunday because it is the weekend, or because she just got over a cold, or her kids are behaving, etc. There are practically unlimited numbers of possible alternative causes, and when the trials are grouped together as in Figure 1b it is more likely that the pattern is merely a coincidence. This example illustrates how the transitions, i.e. the *change* in the cause and effect from one observation to the next convey meaningful information for learning causal strength (e.g., Rottman & Keil, 2012; Soo & Rottman, 2014).

## Learning from States vs. Transitions

Instead of reasoning about transitions, an alternative strategy to learn whether the medicine works is to keep track of the distribution of states experienced. Table 1 summarizes the states in the data from Figure 1a or Figure 1b; both contain the same 7 states just in a different order. By convention, the states are labeled [A], [B], [C] and [D].

Table 1: Frequencies of states in data from Figure 1. Labels for states are shown beside counts in [square brackets].

	Energy = 1	Energy = 0
Drug = 1	2 [A]	1 [B]
Drug = 0	1 [C]	3 [D]

Note. The drug is the cause and energy is the effect.

From Table 1, there are more [A]/[D] states relative to [B]/[C] states, suggesting a positive contingency between the cause and effect (a positive causal relationship). Many models of elemental causal induction compute causal strength from the state frequencies (Table 1). Hattori & Oaksford (2007) documented 41 models of elementary causal induction that use only the state frequencies. These models are intended for "cross-sectional" situations where each observation is independent of the prior one – e.g. observing 7 patients, where three have taken the medicine and four have not. In cross-sectional situations the transitions do not convey meaningful information.

In the current study we are interested in causal learning in longitudinal situations (e.g., tracking one person over time). Because most of the focus within the causal learning

literature has been on cross-sectional rather than longitudinal situations, there is not an existing theory of how people may interpret transitions. In the section below we propose one framework for how learners might interpret transitions and how different types of transitions could influence beliefs about causal strength. We then compare the predictions from this framework to existing models of causal induction. Finally, we present behavioral data from an experiment showing people are sensitive to transitions over and above states, in patterns generally consistent with the transition-based learning framework.

### Transition-based learning

With a binary cause and effect there are four possible states at any given time point. Thus, there are  $4 \times 4 = 16$  possible transitions that can occur between two adjacent time points in a time series. Table 2 categorizes all transitions into different types depending on how consistent they are with a positive causal relationship, a negative one, or no relationship. A transition is consistent with a causal relationship if the transition is likely to be generated by that relationship. In general, with positive causal relationships, changes in the cause (X) are accompanied by changes in the effect (Y) in the same direction (e.g.  $\alpha$  transitions). With negative causal relationships, changes in X lead to changes in Y in the opposite direction (e.g.  $\beta$  transitions). If there is no relationship, changes in X are not associated with changes in Y. From this logic, one can reason backwards to consider how observing a particular transition should influence one's belief concerning the causal relation.

Table 2: Predictions of Transition-Based Learning.

Transitions						Consistent with relation?			$\Delta$
States	Type	$X_0$	$Y_0$	$X_1$	$Y_1$	P+	0	N-	
A to D	$\alpha$	1	1	0	0	✓	✗	✗	++
D to A	$\alpha$	0	0	1	1	✓	✗	✗	++
B to D	$\delta$	1	0	0	0	✓	✓	✗	+
C to A	$\delta$	0	1	1	1	✓	✓	✗	+
B to C	$\beta$	1	0	0	1	✗	✗	✓	--
C to B	$\beta$	0	1	1	0	✗	✗	✓	--
D to B	$\gamma$	0	0	1	0	✗	✓	✓	-
A to C	$\gamma$	1	1	0	1	✗	✓	✓	-
A to B	$\epsilon$	1	1	1	0	✓	✓	✓	0
B to A	$\epsilon$	1	0	1	1	✓	✓	✓	0
C to D	$\epsilon$	0	1	0	0	✓	✓	✓	0
D to C	$\epsilon$	0	0	0	1	✓	✓	✓	0
A to A	$\zeta$	1	1	1	1	✓	✓	✓	0
B to B	$\zeta$	1	0	1	0	✓	✓	✓	0
C to C	$\zeta$	0	1	0	1	✓	✓	✓	0
D to D	$\zeta$	0	0	0	0	✓	✓	✓	0

Note. Each transition is shown to be consistent (✓) or inconsistent (✗) with a positive (P+), negative (N-) or no (0) causal relation.  $\Delta$  is the predicted change in causal strength judgment due to the transition. ++ and -- are large changes to causal strength in the positive vs. negative directions, whereas + and - are smaller changes. 0 is no change to causal strength.

Consider  $\alpha$  transitions – increases in X accompanied by increases in Y ([D to A] transitions), or decreases in X accompanied by decreases in Y ([A to D] transitions). These are transitions that would be generated by a positive causal relationship. Such transitions are unlikely if there were no causal relationship or a negative one – one would need to posit a coincidental hidden cause that influenced Y at the same time that X changed (Figure 1a). Since such transitions are most consistent with a positive relation (not neutral or negative), this framework predicts large positive increases in causal strength judgments after  $\alpha$  transitions.

Next, consider  $\delta$  transitions such as [C to A] – X increases (0 to 1) but Y stays at 1. This transition is consistent with a positive causal relationship with a ceiling effect for Y; it cannot increase any further. A [B to D] transition could be interpreted in the same way but with a floor effect. However, these transitions are also consistent with there being no causal relationship, because a change in X is not accompanied by a change in Y. Because  $\alpha$  transitions are only consistent with a positive causal relationship while  $\delta$  transitions are also consistent with no relation, observing  $\alpha$  should lead to a larger increase in causal strength judgments than observing  $\delta$  (though both should lead to an increase).

$\epsilon$  transitions are when only the effect (Y) changes, while X stays the same. When Y changes, there is no reason to expect X to change regardless of the causal relation. Our framework does not predict change to the causal strength judgment for  $\epsilon$  transitions. In  $\zeta$  transitions, neither X nor Y change. Repeated observations of the same state could be due to the continued causal influence of X, or both X and Y coincidentally remaining in the same state (Figure 1). We predict no change in judgments for  $\zeta$  transitions.

This logic can be extended to transitions consistent with negative causal relationships.  $\beta$  transitions (only consistent with a negative relation) should lead to larger *decreases* than  $\gamma$  transitions (consistent with a negative relation with a floor/ceiling effect, and also with no relation). Observing  $\beta$  should also lead to a larger decrease than observing  $\epsilon$  or  $\zeta$  transitions.

In sum, the most crucial prediction made by this theory is that  $\alpha$  and  $\beta$  transitions will lead to more change (in the positive and negative direction respectively) than the other types of transitions.

### Models of causal induction

We compare the predictions of our framework in Table 2 (the rightmost column) with several existing models of causal strength learning: We briefly present their predictions for longitudinal causal learning in Table 3, with predictions for our transition-based learning (TBL) framework. Many of the models are entirely or largely influenced by the states, so Table 3 groups together the four transitions that end in the same state (gray vs. white). Within each of the four groups in Table 3, the first row are  $\alpha$  or  $\beta$  transitions (both variables change), the second are  $\delta$  or  $\gamma$  (ceiling and floor effects), the third are  $\epsilon$  transitions (effect changes by itself), and the fourth are  $\zeta$  transitions (neither variable changes).

Table 3: Model predictions for 16 transitions.

Transition	$\Delta P$ / PowerPC	RW	TD	TBL	
D to A	++	++	+	$\alpha$	++
C to A	++	++	+	$\delta$	+
B to A	++	++	++	$\epsilon$	0
A to A	++	++	++	$\zeta$	0
A to D	++	0	*	$\alpha$	++
B to D	++	0	*	$\delta$	+
C to D	++	0	0	$\epsilon$	0
D to D	++	0	0	$\zeta$	0
C to B	--	--	-	$\beta$	--
D to B	--	--	-	$\gamma$	-
A to B	--	--	--	$\epsilon$	0
B to B	--	--	--	$\zeta$	0
B to C	--	0	*	$\beta$	--
A to C	--	0	*	$\gamma$	-
D to C	--	0	0	$\epsilon$	0
C to C	--	0	0	$\zeta$	0

Note. ++ and -- denote a predicted increase or decrease that is larger relative to + and - within the same model. 0 denotes no predicted change. \*These cases depend upon too many factors so no generalized predictions can be made.

**$\Delta P$  (Jenkins & Ward, 1965) and Power-PC (Cheng, 1997)** These are two examples of models that are calculated simply from the contingency table (e.g., Table 1). They both produce a causal strength rating from -1 to 1, and can be calculated with the following equations:  $\Delta P = a/(a+b) - c/(c+d)$ , and powPC (for a positive causal strength) =  $\Delta P/[d/(c+d)]$ . If these models are used to calculate causal strength repeatedly after each new observation, then after an A or D observation the causal strength judgment will go up, and after a B or C observation the judgment will go down.

These models are not sensitive to transitions. For example, consider four sequences of data all ending in A: [A,B,C,D,A], [D,A,B,C,A], [C,D,A,B,A], and [B,C,D,A,A]. In all four of these sequences the causal strength ratings from  $\Delta P$  and Power PC would be exactly zero after the 4<sup>th</sup> trial. Then, after the 5<sup>th</sup> trial, the causal strength rating would increase. However, it would increase exactly the same amount under all four sequences. The causal strength would be 1/6 for  $\Delta P$  and 1/3 for Power PC.

**RW (Rescorla & Wagner, 1972; Wagner & Rescorla, 1972)** The Rescorla-Wagner model is a model of associative learning that has also been proposed of human causal learning (Shanks & Dickinson, 1987). RW is a trial-by-trial model of learning that updates weights representing the strength of the association between a cue (cause) and outcome (effect) after each observation. RW was created to model change in associative strength within a single animal over time, so unlike the models mentioned above it is meant to handle time series data.

RW was created to model many temporal phenomena in how associative strengths get updated over time such as acquisition curves and blocking, so unlike the models above it is exquisitely sensitive to the order of the data. Still, RW

works by updating the associative strength after each sequential state observation, [A], [B], [C], or [D]. At any point in time, the change in the associative strength is calculated based on the difference in the error prediction of the outcome (effect) from summing the associative strengths of the present cues (causes). Thus, it does not matter what the immediately previous trial was, the only thing that matters is the current associative strength rating.

One reason why the current associative strength rating matters is that if the current associative strength is zero and an [A] trial occurs, there will be a relatively large increase to the associative strength, but if the current associative strength is .75 and an [A] trial occurs there will be a smaller increase in the associative weight (because the error is smaller). For this reason, we use the prior causal strength rating as an interaction term in all analyses.

Provided that the associative strength is not already at asymptote (+1 or -1), the strength will always increase on an [A] trial and decrease on a [B] trial. Changes to the strength only occur when the cue is present, so no changes occur after [C] and [D] trials; though Van Hamme & Wasserman (1994) have proposed that the strengths be updated even when the cue is absent.

In sum, even though RW is sensitive to many aspects of the order of the observations, the specific prior observation does not have any impact above and beyond the current associative weight. This means that there should be the same amount of change after an [A] trial regardless of the prior state, accounting for the current associative weight.

**TD** Temporal difference (TD) learning is a form of reinforcement learning. Here we discuss a particular instantiation of TD learning that models classical conditioning (Sutton & Barto, 1987). Even though TD is heavily based on RW and it has been widely applied in other areas of psychology (cf. Seymour et al., 2004), as far as we know TD has never been proposed or analyzed as a model of human causal learning. Here we discuss some of the most important differences between RW and TD.

First, whereas RW seeks weights that minimize the prediction of the unconditioned stimulus (effect) at a given instant, TD predicts a sum of future values of the effect signal discounted such that the near future is weighted more than the distant future. Second, learning (or changes to the associative strength) occurs repeatedly moment-to-moment within a trial as opposed to just once at the end of the trial. Third, whereas RW only updates weights for the cause when the cause is present, TD updates weights for the cause in proportion to the strength of an “eligibility” trace (similar to a memory/salience trace) of the cause. If the cause has been present for a while, learning is fast. But if it was recently absent, learning is slow until it is more eligible. Even after the cause disappears some learning can occur to the extent that the trace persists. Fourth, the weights for the cues are not bounded; we just focus on whether the weights change in the positive or negative direction.

The dynamics of all of these features plus others means that (unlike the other models discussed so far) TD actually

makes predictions about transitions between specific states (Table 3). First, TD predicts a greater increase for [A to A] and [B to A] than [C to A] or [D to A] transitions. In the latter two transitions the trace of the cause is initially zero, and it takes time to become activated, slowing down learning. In the former two transitions the cause is already present from the previous trial so the eligibility trace is initially higher, speeding up learning. The difference between [B to B] and [A to B] vs. [C to B] and [D to B] is also due to the eligibility trace, just in the negative direction.

In the transitions [C to C], [C to D], [D to C], and [D to D], the cause is never present so its eligibility is always zero and its weight is not updated. This is similar to how RW does not update strength when the cue is absent.

The transitions [A to B], [D to B], [B to C], and [A to C] are all extremely dynamic and depend on the prior weight (above vs. below zero) and the prior weight of the unobserved cue (above vs. below zero). Because of the extreme level of the dynamics we cannot make a generalized characterization of how the weights get updated for these transitions.

**Comparisons Between Models** Predictions for the transition-based learning (TBL) framework are included in the right column in Table 3 (compare to Table 2). There are several comparisons that shed light on the similarities and differences between the models.

First, there is some consistency in the models. Transitions ending in A and D are viewed as positive (or neutral) evidence for all models, whereas those ending in B and C are negative (or neutral) evidence for all models. Second,  $\Delta P$ , Power PC and RW make the same predictions for all transitions ending in the same state. In contrast, TD and TBL make different predictions for transitions that end in the same state. Third, even though both TD and TBL are sensitive to transitions, the predictions are nearly opposites. Consider the transitions ending in A. TD predicts larger increases for [B to A] and [A to A] than [D to A] and [C to A]. TBL makes the exact opposite predictions. Most importantly, TBL predicts the largest increase for [D to A]. This same basic pattern also plays out in the transitions ending in B. Here we are ignoring the transitions ending in C and D because TD does not make general predictions. The goal of our experiment was to test which of these models predicts the trial-by-trial changes in causal strength judgments the best.

## Experiment

Subjects observed sets of longitudinal data and made causal strength judgments after each trial. We were focused on whether the changes in causal strength judgments from trial to trial were influenced by the immediately-prior trial.

### Methods

**Subjects** 100 subjects were recruited through Amazon Mechanical Turk (MTurk) and paid \$1.75 for completing the experiment. The experiment was conducted online and took roughly 15-20 minutes to complete. Because some

subjects began the study but stopped midway, we actually collected data from 177 subjects though some of their data was partial (i.e. less than the full number of scenarios).

**Design and stimuli** Subjects were presented with sets of data consisting of a binary cause and effect. Each subject viewed 15 sets of data, and each set (one scenario) consisted of 8 trials. For each participant five data sets had  $\Delta P = 0$  (2 observations of each state), another five had  $\Delta P = -0.5$  (1 observation each of [A] and [D], 3 observations each of [B] and [C]), and another five had  $\Delta P = 0.5$  (3 observations each of [A] and [D], 1 observation each of [B] and [C]). The trials within a data set were randomly ordered. The reason for having data sets with positive, negative, and neutral contingencies was to have a sampling of all the transitions. For example,  $\beta$  transitions are very rare in the  $\Delta P = 0.5$  data sets. Each data set had 8 states, and thus 7 transitions.

**Procedure** Subjects were told to imagine they were researchers studying the effects of drugs on chemicals in the blood of monkeys. Each scenario involved testing one drug on one chemical in one monkey, with a new drug, chemical, and monkey for each of the 15 scenarios.

Each trial within a scenario involved the drug either being administered or not (using an intravenous drip), and then participants were told that one hour later a blood test is conducted to reveal whether the chemical is high or low. After seeing the blood test participants estimated the causal strength of the drug on the target chemical using a slider with the following anchors: -99 = ‘When the drug is on, the chemical is usually low. When the drug is off, the chemical is usually high’, 0 = ‘There is no relationship between whether the drug is on or off and the level of the chemical’, and 99 = ‘When the drug is on, the chemical is usually high. When the drug is off, the chemical is usually low’. The slider stayed at the same value from the participant’s judgment after the prior trial; after each trial participants could either move the slider to update their judgment or a check-box to keep the same judgment from before. After making the judgment participant saw the next trial (whether the drug was administered or not, and the blood test) until they were finished with the 8 trials. Subjects completed a practice scenario and 15 actual scenarios (5 from each contingency); the order of the 15 scenarios was random.

### Results

36 subjects were excluded from the analysis because their responses indicated a misinterpretation of the scale – on [D] states; they always reduced their causal strength judgments. Further investigation revealed that these subjects’ judgments tracked the occurrence of the effect; they increased judgments when the effect was present (even on [C] trials), and decreased them when it was absent. This interpretation occurred despite our best efforts at defining a positive vs. negative relation (see methods section). We did not want to further train subjects on the use of the scale by providing feedback because we did not want to imply that there was one right answer and wanted to preserve their natural use of the scale as much as possible. However, this interpretation

of the scale is not explainable by any of the models of causal induction. Thus, we eliminated from the analysis subjects who increased at least half their judgments on [A to C] and [B to C] transitions *and* decreased at least half of them on [A to D] and [B to D] transitions.

We analyzed the change in causal strength judgments (the difference between the judgment at the present and prior trials) for each of the 16 types of transitions. We ran four regressions, one for each group of transitions ending in the same state (shaded rows in Table 4). The  $\alpha$  and  $\beta$  transitions (top row within each of the 4 sets in Table 4) were treated as the reference transition, because the most important hypothesis was whether the  $\alpha$  and  $\beta$  transitions produced larger changes than the other transition types. A by-subject random intercept and a random slope on transition type were included to account for the fact that each individual made multiple judgments for each transition type.

One challenge in analyzing change scores is that the causal strength judgment at the prior trial can constrain the amount of change. For example, for transitions expected to lead to an increase, a very high starting point would constrain the amount of possible increase. It is also possible that different transition types would produce different amounts of change at different prior strength levels, so an interaction between prior strength and transition type was included in the regression. (These interactions are not discussed further in the current manuscript.)

Table 4 displays the results of the four regressions. Within the four transitions ending in a given state, the other three transitions were all compared against the top transition. In typical regression tables the difference between levels is reported. But for ease of interpretation we translate the differences to their own group means (e.g. [C to A] transitions produced an average increase of 14 points). The rightmost column in Table 4 summarizes which of the three transitions are significantly different in size compared to the top transition (+ is a significant smaller increase than ++).

The first impression of this table is that all the transitions ending in A and D produced increases in causal strength, whereas those ending in B and C produced decreases. This finding is most consistent with Power PC and  $\Delta P$ .

The second striking finding is that the  $\alpha$  and  $\beta$  transitions always produced more extreme changes compared to  $\zeta$  transitions (forth row of each set). This finding is uniquely predicted by TBL; it is not predicted by the other models and is even the opposite prediction made by TD.

Third, the other predictions made by TBL, that  $\alpha$  and  $\beta$  would be stronger than  $\delta$  and  $\gamma$  (ceiling and floor effects), and  $\varepsilon$  (when the effect changes by itself) were not supported (except one instance of  $\delta$ ). The following paragraphs go through the results in more detail.

Amongst transitions ending in [A], all four transitions lead to increases. [D to A] transitions led to the largest increase, an average increase of 21. Compared to [D to A], [C to A] transitions and [A to A] transitions led to significantly smaller increases. The significance can be seen by examining whether the 95% CI includes the mean change

for [D to A]. For example, 95% CI for [C to A], [8, 19] is entirely lower than the mean for [D to A], 21. [A to A] also produced a significantly smaller increase than [D to A], but [B to A] was not significantly different.

Among the transitions ending in [D], the  $\alpha$  transition was significantly stronger than  $\zeta$ . It was not significantly different than  $\delta$ , and it actually was weaker than  $\varepsilon$  (which is not predicted by any model).

Among the transitions ending in [B] and [C], the  $\beta$  transitions were more extreme than the  $\zeta$  transitions, but the  $\beta$  transitions were the same size as the  $\gamma$  and  $\varepsilon$  transitions.

In summary, the results show some patterns consistent with Power PC and  $\Delta P$ , as well as one consistent pattern uniquely predicted by TBL, that  $\alpha$  and  $\beta$  transitions resulted in larger changes to causal strength than  $\zeta$  transitions.

Table 4: Regression results for effect of transitions. Separate models for each group of transitions by end-state.

Transition Type	Mean Change	95% CI of Change		Summary
		Lower	Upper	
D to A $\alpha$	21	18	24	++
C to A $\delta$	14	8	19	+
B to A $\varepsilon$	19	13	25	++
A to A $\zeta$	10	7	14	+
A to D $\alpha$	11	8	14	++
B to D $\delta$	10	4	15	++
C to D $\varepsilon$	17	11	24	+++
D to D $\zeta$	7	3	9	+
C to B $\beta$	-16	-19	-14	--
D to B $\gamma$	-15	-21	-9	--
A to B $\varepsilon$	-19	-25	-12	--
B to B $\zeta$	-12	-15	-8	-
B to C $\beta$	-15	-18	-13	--
A to C $\gamma$	-13	-18	-7	--
D to C $\varepsilon$	-14	-21	-10	--
C to C $\zeta$	-11	-14	-7	-

Note. The Change column indicates the relative increase or decrease of a particular transition relative to the top row in the same group of transitions.

## General Discussion

Previous work studying elemental causal learning has focused on how causal strength is learned from states – [A], [B], [C] and [D]. In the current article we proposed an extreme version of an elemental causal learning theory that focuses exclusively on transitions. The proposal we put forth was intended to be provocative – to theorize how different transitions could be interpreted completely independently of states. In reality, we are not proposing that people exclusively rely upon transitions and indeed the results suggest a combination of strategies.

The main finding in support of the transition-based learning theory was that when both the cause and effect changed ( $\alpha$  and  $\beta$ ) people changed their causal strength judgments more than when the same state was repeated ( $\zeta$ ), controlling for the prior causal strength judgment. Another

way to think about these findings is that when the same state is repeated ( $\zeta$ ) there actually is no “transition”; the repeat of the state could be viewed as redundant (i.e. repeated observations of a state can be collapsed into a single extended data point). This finding was the basic phenomenon from the example involving Figure 1.

The TBL framework made two predictions that were not supported. We review these predictions because they reveal some interesting reasoning habits. The first prediction involved  $\delta$  and  $\gamma$  transitions. One example is when the cause changes from 0 to 1 but the effect stays at 1. We predicted that this transition could be viewed as consistent with a positive causal relation (a ceiling effect), however it could also be consistent with no relation – the cause turns on but since the effect is already on the cause was not responsible for the effect. Out of the four  $\delta$  and  $\gamma$  transitions one was weaker than the  $\alpha$  and  $\beta$  transitions, but the other three were not significantly different. This suggests that the participants were attributing the state of the effect to the cause even though the state of the effect was present before the cause.

The second disconfirmed prediction made by TBL was that transitions when the effect changes on its own would not result in changes to the causal strength (e.g., if the cause stays at 1, and the effect changes from 0 to 1). The reasoning was that it is always possible that the effect could change on its own due to some unobserved factor, but this change should be attributed to the unobserved factor, not the target cause. Another interpretation is that there really is a positive causal relation, but at the initial state there was a temporary unobserved inhibitory factor.

From the discussion above, it is evident that transitions can be interpreted in multiple ways. There is some interesting research on how observed states can be interpreted differently given different prior knowledge – in some instances even [A] can be interpreted as negative evidence (Luhmann & Ahn, 2011). When reasoning about states, beliefs about unobserved factors drive the different interpretations. And as seen in the paragraphs above, we also hypothesize that beliefs about unobserved factors could be responsible for different interpretations. More fully developing a transition-based theory of causal induction will require clarifying the interpretations of the transitions, which could be facilitated by eliciting verbal explanations of the interpretation of a given transition.

The current results suggest that elemental causal learning in longitudinal contexts involves a combination of transition and state-based reasoning. One important goal for future research is to better capture how these two types of reasoning get used – do they get used simultaneously, are there individual differences, or does a single learner sometimes focus on one interpretation and other times focus on another? Future research will investigate the factors that promote the use of one reasoning pattern over another. This will ultimately result in a more complete theory of real-world elemental causal learning that (unlike most existing theories) makes the distinction between data in cross-sectional and longitudinal contexts.

## References

- Cheng, P. W. (1997). From covariation to causation: A causal power theory. *Psychological Review*, *104*(2), 367–405. doi:10.1037//0033-295X.104.2.367
- Hattori, M., & Oaksford, M. (2007). Adaptive non-interventional heuristics for covariation detection in causal induction: model comparison and rational analysis. *Cognitive Science*, *31*(5), 765–814. doi:10.1080/03640210701530755
- Jenkins, H. M., & Ward, W. C. (1965). Judgment of Contingency Between Responses and Outcomes. *Psychological Monographs: General and Applied*, *79*(1), 1–17.
- Luhmann, C. C., & Ahn, W.-K. (2011). Expectations and interpretations during causal learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *37*(3), 568–87. doi:10.1037/a0022970
- Rescorla, R. A., & Wagner, A. R. (1972). A Theory of Pavlovian Conditioning: Variations in the Effectiveness of Reinforcement and Nonreinforcement. In *Classical Conditioning II: Current Theory and Research* (pp. 64–99).
- Rottman, B. M., & Keil, F. C. (2012). Causal structure learning over time: observations and interventions. *Cognitive Psychology*, *64*(1-2), 93–125. doi:10.1016/j.cogpsych.2011.10.003
- Seymour, B., Doherty, J. P. O., Dayan, P., Koltzenburg, M., Jones, A. K., Dolan, R. J., ... Frackowiak, R. S. (2004). Temporal difference models describe higher-order learning in humans. *Nature*, *429*(June), 664–667. doi:10.1038/nature02636.1.
- Shanks, D. R., & Dickinson, A. (1987). Associative accounts of causality judgment. In G. H. Bower (Ed.), *Psychology of learning and motivation: Advances in research and theory* (pp.229-261). San Diego, CA: Academic Press.
- Soo, K. W., & Rottman, B. M. (2014). Learning Causal Direction From Transitions With Continuous And Noisy Variables. In P. Bello, M. Guarin, M. McShane, & B. Scassellati (Eds.), *Proceedings of the 36th Annual Conference of the Cognitive Science Society* (pp. 1485-1490).
- Sutton, R. S., & Barto, A. G. (1987). A Temporal-Difference Model of Classical Conditioning. In *Proceedings of the 9th Annual Conference of the Cognitive Science Society*.
- Van Hamme, L. J., & Wasserman, E. A. (1994). Cue Competition in Causality Judgments: The Role of Nonpresentation of Compound Stimulus Elements. *Learning and Motivation*, *25*, 127–151. doi:10.1006/lmot.1994.1008
- Wagner, A. R., & Rescorla, R. A. (1972). Inhibition in Pavlovian Conditioning: Application of a Theory. In *Inhibition and Learning* (pp. 301–336).