# Causal learning from interventions and dynamics in continuous time

**Neil R. Bramley**[1] (neil.bramley@nyu.edu), **Ralf Mayrhofer**[2] (rmayrho@gwdg.de)
**Tobias Gerstenberg**[3] (tger@mit.edu), **David A. Lagnado**[4] (d.lagnado@ucl.ac.uk)

[1]Department of Psychology, NYU, New York, NY, 10003, USA
[2]Department of Psychology, University of Göttingen, Gosslerstr. 14, 37073, Germany
[3]Department of Brain and Cognitive Sciences, MIT, Cambridge, MA 02139, USA
[4]Department of Experimental Psychology, UCL, London, WC1H 0DS, UK

## Abstract

Event timing and interventions are important and intertwined cues to causal structure, yet they have typically been studied separately. We bring them together for the first time in an experiment where participants learn causal structure by performing interventions in continuous time. We contrast learning in acyclic and cyclic devices, with reliable and unreliable cause–effect delays. We show that successful learners use interventions to structure and simplify their interactions with the devices and that we can capture judgment patterns with heuristics based on online construction and testing of a single structural hypothesis.

**Keywords:** causal learning; intervention; time; causal cycles; structure induction; dynamics.

In a dynamically unfolding world, using actions to uncover causal relationships requires good timing. It is hard to tell whether a new medication is effective if you take it with others, or just as you start to feel better. Likewise, it is hard to tell whether a new law lowers crime if it is introduced just after other reforms or before a major election. Such inferences, having to do with delayed effects and an evolving causal background, can be particularly tough in cyclic systems in which feedback loops make prediction difficult even with complete knowledge (Brehmer, 1992). Thus, for interventions to be effective tools for unearthing causal structure it is important to time and locate them carefully, paying close attention to the temporal dynamics of surrounding events and the possibility of feedback loops.

Previous work has shown that people make systematic use of temporal information, taking event order as a strong cue to causal order (Bramley, Gerstenberg, & Lagnado, 2014), and making stronger attributions when putative cause–effect delays are in line with expectations (Buehner & McGregor, 2006) and have low variance across instances (Greville & Buehner, 2010). Recent work has also developed frameworks for probabilistic causal inference from event timings based on parametric assumptions about cause–effect delays (Bramley, Gerstenberg, Mayrhofer, & Lagnado, submitted; Pacer & Griffiths, 2015).

A distinct line of work has shown that people are adept at inferring causal structure from interventions — idealized actions that set variables in a system (e.g., Bramley, Dayan, Griffiths, & Lagnado, 2017; Coenen, Rehder, & Gureckis, 2015). This work has not explored the role of temporal information however. While researchers have speculated about the close relationship between temporal and interventional inference (e.g., Lagnado & Sloman, 2004), our paper is the first to explore interventional causal learning in continuous time.

## The learning problem

We explore the general problem of how people learn about a causal system by interacting with it in continuous time. We focus on abstract causal "devices" made up of 3–4 components (cf. Figure 1). For causally related components, we assume each activation of a cause will tend to bring about a single subsequent activation of its effect after a parametric delay (described below). For example, Figure 1a shows a learner's interactions with a $B \leftarrow A \rightarrow C$ Fork during which time they perform four interventions. Activations of both $B$ and $C$ succeed the interventions on $A$ but with some variability in delays.

We focus on situations where components never spontaneously activate, but where causal relations work stochastically (e.g., are successful with probability $w_S$). Any pair of components can be connected in either, neither or both directions resulting in a hypothesis space $S$ of 64 possible structures for devices made up of three components, and 4096 for four components. Learners can intervene on the devices by directly activating any component at any moment of their choosing. Interventions are always successful in that they instantaneously activate the targeted component. The downstream causal effects of intervened-on components are the same as those of components that were activated by other components. Thus, we model the consequences of interventions in analogy to the Do(.) operator introduced by Pearl (2000), such that interventions provide no information about the causes of the intervened-on component.

### Choosing interventions

Seeing the effects of one's interventions in continuous time provides rich information for causal inference. On the flip side, there are also no completely independent trials. For instance, in Figure 1a, the early interventions on $C$ and $B$ might, in principle, be responsible for the observed effects that happen shortly after the intervention on $A$. In general, one cannot rule out the possibility something that happened earlier is still exerting its influence, or that an effect is yet to reveal itself. Fortunately, interventions provide anchor points. We know that events due to interventions weren't caused by anything else, and that these events only affect the future but not the past (Lagnado & Sloman, 2004). This means that by intervening, learners can recreate some of the advantages that come with a discrete trial structure. For example, by waiting long enough between interventions to be confident prior effects have dissipated, an otherwise confusing event stream
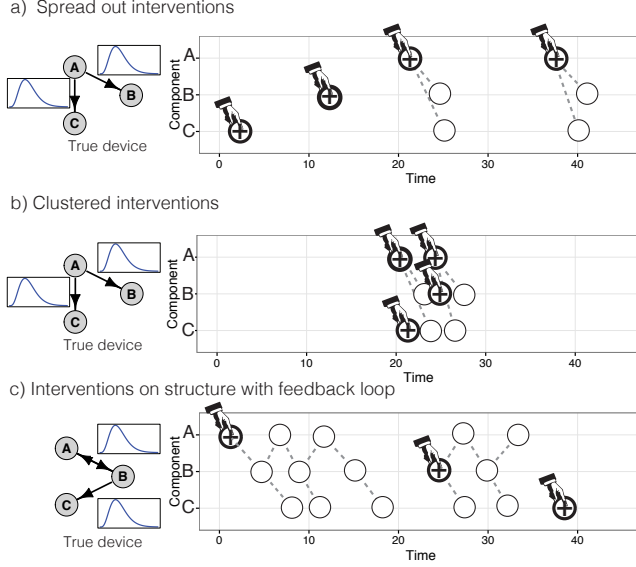
Figure 1: Examples of using real-time interventions to infer causal structure. Left: True generative causal model with subplots showing delay distributions. Right: Timelines showing an active learners' interactions with each system with a row for each component *A* (top), *B* (middle) and *C* (bottom), and white circles indicating their activations over 45 seconds (x-axis). "+" symbol and incoming hand icon indicate interventions. Dashed gray lines indicate the actual cause–effect relationships.

becomes more palatable and informative about the underlying structure. Figure 1b gives an example of interventions that are not well chosen. The learner performs four interventions in the same locations as Figure 1a but does so in close succession. It is hard to attribute causal responsibility for these activations, since there are so many similarly plausible candidates. Consequentially, this data is considerably less informative.

In discrete-trial interventional learning, participants exhibit a *positive testing strategy* — they prefer to intervene on root variables that bring about many effects (Coenen et al., 2015). While often not leading to the most globally informative choice, a positive testing strategy is an effective way of assessing the adequacy of one's current working hypothesis, making it a manifestation of confirmatory testing (Nickerson, 1998). Many other components will be affected if one's hypothesis is right, and few if it is wrong. Repeated positive testing might be more justifiable in the continuous time context because cause–effect delays may play out differently each time, and potential temporal reversals between variable activations will help to rule out candidate structures (Bramley et al., 2014). For example, in Figure 1a the second intervention on *A* leads to *B* and *C* occurring in reversed order, allowing the learner to rule out a $A \rightarrow B \rightarrow C$ Chain structure.

## Causal cycles

The vast majority of causal learning studies have focused on acyclic causal systems in which causal influences flow only in one direction, never revisiting the same component. However, many natural processes are cyclic and people frequently report cyclic relationships when allowed to do so (e.g.

Sloman, Love, & Ahn, 1998). While there are ways of adapting the causal Bayes net formalism to capture cycles (Rehder, 2016), these generally simplify the problem to influences between fixed time steps (e.g. Rottman & Keil, 2012), or just to the long-run equilibrium distribution (e.g. Lauritzen & Richardson, 2002). However, by focusing on continuous time and developing a representation capable of modeling causal dynamics, we are able to directly compare learning in acyclic and cyclic causal systems.

Dynamic systems can be hard to predict even with perfect knowledge. Positive feedback loops can lead to sensitive dependence on initial conditions with very different behavior resulting from small perturbations in starting conditions (e.g., Gleick, 1997). Figure 1c gives an example of interventions on a cyclic causal system (assuming that the connections work 90% of the time). Interventions initialize looping behavior because of the bidirectional relationship $A \leftrightarrow B$ (e.g., $A \rightarrow B \rightarrow A \rightarrow B \ldots$) leading to many subsequent activations of both the loop components and the output component *C*, continuing until either the $A \rightarrow B$ or $B \rightarrow A$ connection fails. Based on simply looking at the timeline, it seems likely that it will be easier to identify which components are either directly involved in cycles, or outputs from cyclic components (due to their recurrent activations), but harder to identify the exact causal relationships (e.g. whether it is *A* or *C* that causes *B* in this example since both tend to recur shortly before *B*).

## Normative inference

As a benchmark, we developed a Bayesian model of causal structure inference. We consider the data $\mathbf{d}_\tau \left\{ d_X^{(1)}, \ldots, d_X^{(n)} \right\}$ to be made up of all activations (with events indexed in chronological order and *X* indicating the activated component) conditioned upon the set of interventions $\mathbf{i}_\tau = \left\{ i_X^{(1)}, \ldots, i_X^{(m)} \right\}$. Both $\mathbf{d}_\tau$ and $\mathbf{i}_\tau$ are restricted to the interval between the beginning of the clip and time $\tau$, which we assume to be the moment at which the learner makes the inference. For instance, one might interact with a causal device for 5000 ms, performing interventions on components *A* and *B* at 100 ms and 1200 ms respectively: $\mathbf{i}_{5000} = \{ i_A^{(1)} = 100, i_B^{(2)} = 1200 \}$, and observing two activations of *C*: $\mathbf{d}_{5000} = \{ d_C^{(1)} = 1500, d_C^{(2)} = 2800 \}$.

Normative Bayesian structure inference involves updating a prior over structure hypotheses $P(S)$ with the likelihood $p(\mathbf{d}_\tau | S; \mathbf{i}_\tau, \mathbf{w})$ to get a posterior belief over structures $P(S | \mathbf{d}_\tau; \mathbf{i}_\tau, \mathbf{w})$ given the set of parameters $\mathbf{w}$:[1]

$$P(S | \mathbf{d}_\tau; \mathbf{i}_\tau, \mathbf{w}) \propto p(\mathbf{d}_\tau | S; \mathbf{i}_\tau, \mathbf{w}) \cdot P(S) \quad (1)$$

An immediate issue with calculating the likelihood of an observed set of activations given a candidate model is that there are likely to be multiple potential paths of *actual causation* that could have produced the data (Halpern, 2016), each

---

[1] In this specific case, we assume the parameters (i.e., causal strength $w_S$, expected length of delays $\mu$, and delay variability $\alpha$) to be known which is consistent with the setup of the experiment.
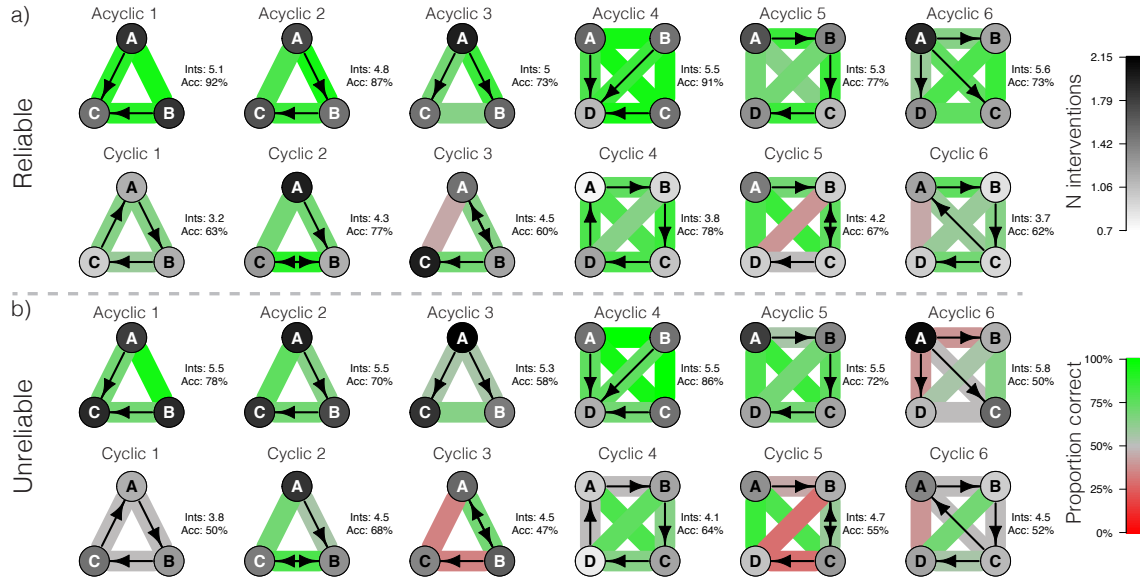
Figure 2: Devices tested and results from experiment in a) *reliable* and b) *unreliable* delay conditions. **Node shading:** Intervention choice prevalence by component. **Edge shading:** accuracy. *Note*: Ints = average number of interventions performed; Acc = mean accuracy.

of which implying a different likelihood. For example, if the true structure is a $A \to C \gets B$ Collider, the data above might be produced in two ways. $A$ could have caused the first activation of $C$ and $B$ the later ($i_A^{(1)} \to d_C^{(1)}, i_B^{(1)} \to d_C^{(2)}$). Alternatively, $A$ could have caused the later activation of $C$ and $B$ the earlier ($i_A^{(1)} \to d_C^{(2)}, i_B^{(1)} \to d_C^{(1)}$).

However, as there can only be one true path of actual causation in the set of possible paths $\mathbf{Z}_s$, we can sum over these to get the likelihood of the data given a candidate model $s \in S$:

$$p(\mathbf{d}_\tau | s; \mathbf{i}_\tau, \mathbf{w}) = \sum_{\mathbf{z}' \in \mathbf{Z}_s} p(\mathbf{d}_\tau | \mathbf{z}'; \mathbf{i}_\tau, \mathbf{w}) \qquad (2)$$

We assume that the actual causal delays (in $\mathbf{Z}_s$) are Gamma distributed (see also Bramley et al., submitted) with a known expected duration $\mu$ and shape $\alpha$ (i.e., variability). The likelihood of the data given a specific path $\mathbf{z}'$, then, is the product of the (Gamma) likelihoods of the observed delays and causal strength $w_S$ combined with the likelihoods of (non-)events, the occurrence of which failed either due to the $1 - w_S$ causal failure rate or due to the effect potentially occurring after $\tau$ (i.e., some time in the future).

With these ingredients the posterior belief over causal structure hypotheses can be determined. However, it is only feasible to enumerate all possible paths of actual causation for a sufficiently small number of events. While for a large number of events the calculations become intractable, we were able to compute the posteriors in the described manner for the data from the current experiment, resorting only in rare cases to an approximation.[2]

## Experiment

Participants' task was to discover the causal connections between the components of several devices in limited time

(see Figure 2). Half of the devices were *acyclic* (top; no feedback loops) and half were *cyclic* (bottom; contained a feedback loop). Participants were able to activate any of the components by clicking on them. We were interested in how participants chose *where* to intervene and *when*. We examined two delay conditions between subjects, one in which the true cause–effect delays were *reliable* (Gamma distributed with $\alpha = 200, M \pm SD$ $1.5 \pm 0.1$ seconds) and one where they were *unreliable* ($\alpha = 5, M \pm SD$ $1.5 \pm 0.7$ seconds). Following Greville and Buehner (2010), we expected that performance would be better when causal delays were *reliable*. We also predicted that complex dynamics would lead to worse performance when the true structure was *cyclic*, and that successful participants would spread their interventions widely over time, thus minimizing the ambiguity of resulting patterns of effects.

## Methods

**Participants** Forty participants (14 female, aged $32 \pm 9.0$) were recruited from Amazon Mechanical Turk (yielding 20 subjects in each delay-reliability condition) and were paid between \$0.50 and \$3.20 (\$2.06 $\pm$ 0.39) depending on performance (see Methods section). The task took around 20 minutes.

**Materials and procedure** Each device was represented with a circle for each component and boxes marking the locations of the potential connections (see Figure 3a).[3] Trials lasted for 45 seconds during which components activated if clicked on or if caused by the activation of another component, with delay and probability governed by the true underlying network (Figure 3b). Causal relationships worked 90% of the time (i.e., causal strength $w_S = 0.9$) and there were no spontaneous activations. Activated components turned yellow for 200ms, and intervened-on components were additionally marked by a "+" symbol. Initially, all components were inactive and no

---

[2]Where necessary, we ruled out paths that implied an implausibly high number of failed connections, or extreme cause–effect delays, until the number of possible paths fell below 100,000.

[3]Try the task `https://www.ucl.ac.uk/lagnado-lab/el/it` or watch a trial `https://www.ucl.ac.uk/lagnado-lab/el/itv`.
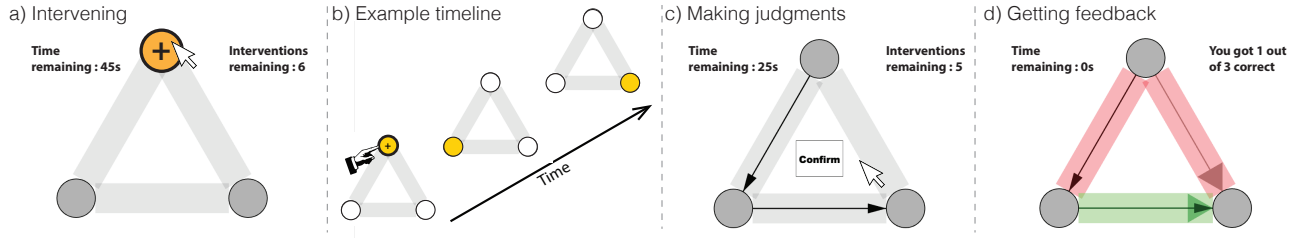
Figure 3: Experimental procedure. a) Up to 6 interventions could be performed by clicking on the components during the 45 second trial. b) This would lead to subsequent activations determined by causal connections and delays in the true model. c) Participants marked their beliefs about the structure during the trials by clicking on the edges. d) At the end of each trial they received feedback. Broad gray arrows: ground truth, Green = correct, Red = incorrect.

connections were marked between them.

Prior to the inference tasks, participants were trained on the delays in their condition and how to register structure judgments through interaction with an an example device. They then had to correctly answer comprehension check questions and complete a practice problem, before facing the 12 test devices in random order with randomly orientated and unlabeled components.

In the test phase, participants could perform up to 6 interventions on each trial and register/update their judgments about the causal structure as often as they liked until the 45 seconds for a device ran out (for details see Figure 3). At the end of each trial, they were given feedback showing the true relationships and which of them they had correctly identified. To incentivize proper judgments, bonuses were paid based on connections participants had registered at a randomly chosen point during each trial.

### Results

We analyze participants' judgments by first comparing their accuracy by delay-reliability condition (between subjects: *reliable* vs. *unreliable*) and device type (within subject: *acyclic* vs. *cyclic*). We then analyze the timing and spacing of participants' interventions and how these relate to the evidence and judgments.

**Accuracy** Participants updated and confirmed their judgment about the structure M±SD $1.6 \pm 1.2$ times per trial on average. Judgment time was not significantly related to accuracy, but within trials, final judgments were slightly more accurate than initial judgments, with participants correctly identifying $69\% \pm 30\%$ (chance performance would be 25%) compared to $65\% \pm 28\%$ of the connections, $t(479) = 5.2, p < .001$ (remember that bonuses incentivised making judgments early). Only 4% of judgment updates decreased the number of connections, 24% resulting in the same number as before, and 72% increasing the number of connections.

Focusing on final judgments, participants correctly identified [*reliable,acyclic*]: $82\% \pm 29\%$, [*reliable,cyclic*]: $68\% \pm 28\%$, [*unreliable,cyclic*]: $69\% \pm 29\%$, [*unreliable,cyclic*]: $56\% \pm 29\%$ of the connections. A repeated measures analysis revealed a significant effect of delay-reliability condition, $F(1, 38) = 4.6, p = .04$, and cyclicity, $F(1, 38) = 39, p < .001$, but no interaction, with *unreliable* delays and *cyclic* structures associated with lower accuracy. Figure 2 shows that participants found the Cyclic 3, 5 and 6 structures hardest to identify on average, struggling in particular with distin-

guishing looping from output components.

Ideal Bayesian inference based on the evidence generated by participants predicts a different pattern. While *reliable* delays allow greater accuracy than *unreliable* ones, $F(1, 38) = 24.3, p < .001$, there is no predicted difference in accuracy between *acyclic* and *cyclic* devices, $F(1, 38) = 0.43, p = .5$. In fact, posterior uncertainty over all possible models, measured by Shannon entropy, was generally lower for evidence generated by a *cyclic* $.74 \pm 1.26$ than an *acyclic* $1.95 \pm 1.29$ devices, $F(1, 38) = 109, p < .001$.

**Timing of interventions** We hypothesized that spacing interventions out in time would be important for successful learning. Participants waited $7.3 \pm 2.8$ seconds between interventions on average. In a regression including delay condition and total number of interventions as covariates, leaving longer intervals between interventions was positively associated with accuracy, $F(1, 36) = 14.0, \beta = 0.04, \eta_p^2 = .26, p = .001$, with no interaction with condition. The variability of these gaps — measured by their coefficient of variation $CV = \frac{\sigma}{\mu}$ — was also inversely related to accuracy, $F(1, 36) = 7.9, \beta = -0.5, \eta_p^2 = .18, p = .008$ and this effect was stronger in the *unreliable* delay condition, $F(1, 35) = 4.5, \eta_p^2 = .11, p = .04$. We also assessed the intervals participants left after the most recently preceding event (whether this was an intervention or an effect) before performing their next intervention. Again larger intervals, $F(1, 36) = 7.7, \beta = 0.06, \eta_p^2 = .18, p = .008$, and less variation, $\beta = -.25, F(1, 36) = 5.0, \eta_p^2 = .12, p = .03$, was associated with accuracy with neither measure interacting with delay condition. Both larger intervals between interventions, and between interventions and the most recently preceding effect were also associated with lower posterior entropy, with $\beta = 0.05, F(1, 36) = 9.9, \eta_p^2 = .22, p = 0.003$ and $\beta = 0.09, F(1, 36) = 8.1, \eta_p^2 = .18, p = 0.007$, respectively. However, there was no evidence for a relationship between entropy and the variability of either interval type.

**Positive testing** We found evidence of a preference for positive testing, with participants performing $1.2 \pm 0.5$ times as many interventions per root component than per non-root component $t(59) = 3.9, p < .001$. This preference was associated with higher accuracy after accounting for condition, $F(1, 37) = 21, \eta_p^2 = 0.37, p < .001$, and did not interact with condition. Degree of root preference, however, was not significantly related to posterior uncertainty from the perspective of an ideal Bayesian learner.

**Adaptation to cycles** While participants performed fewer interventions on *cyclic* ($4.1 \pm 1.1$) compared to *acyclic* ($5.4 \pm 0.7$) devices, $t(39) = 8.7, p < .001$ (see Figure 2), they still experienced far more effects in the *cyclic* systems ($29.3 \pm 10$) compared to the *acyclic* ones ($4.7 \pm 1.1$), $t(39) = 15.5, p < .001$. This was due to the reciprocal relationships sustaining activations until one of the links failed. Thus while there was normatively more evidence available in the cyclic trials — as reflected by the generally lower posterior uncertainty — the large number of events resulted in more ambiguous evidence, with many candidate causes per effect and a large number of potential actual causal pathways.

**Summary** Participants were better at identifying causal relations from interventions when delays were *reliable* and the true structure was *acyclic*. Meanwhile, ideal learner accuracy was affected by reliability by not cyclicity. Successful participants spread their interventions out more in time, waited longer after previous events, distributed them more evenly and favored root components. Participants frequently updated their models by adding additional connections but rarely removed connections.

## Modeling heuristic inferences

Participants' deviations from the prediction of an ideal Bayesian learner suggests that they relied on simpler learning strategies. In this section we compare judgment patterns to several heuristic models inspired by work on order–driven (e.g., Bramley et al., 2014) and incremental causal structure learning (e.g., Bonawitz, Denison, Gopnik, & Griffiths, 2014; Bramley et al., 2017).

Several papers have proposed that human causal learning is based on the adaptation of a single global hypothesis (Bonawitz et al., 2014), which might be achieved incrementally through making local changes as data is observed (Bramley et al., 2017). This seems particularly applicable in a continuous-time context, where normative inference is tough and the evidence arrives continuously. People may learn locally, ignoring dependence on beliefs about surrounding relationships (e.g. Fernbach & Sloman, 2009), or use their current model as a basis, comparing observations against predictions, only adding new connections to explain events that cannot easily be accommodated by their existing model (Bramley et al., 2017).

The idea that learners might construct their causal hypotheses incrementally can be combined with different degrees of sensitivity to timing as well as the predictions of their current structure hypothesis. This suggests several potential heuristics that adapt a single model belief $b$ as events are experienced. The result in each case is a single structural belief that evolves as events occur (we write $\mathbf{b} = \{b^{(0)}, \ldots, b^{(n)}\}$, where the sequence of belief indices correspond to the event indices in $\mathbf{d}_\tau$):

1. **Order Only (OO)** Heuristic OO attributes each new effect to the most recently preceding event at any different component (either the most recent intervention in $\mathbf{i}_\tau$ or activation in $\mathbf{d}_\tau$). If the currently held model hypothesis $b^{(t-1)}$ does not contain a respective edge, $b^{(t-1)}$ is augmented with an edge to make $b^{(t)}$. Figure 4a gives an example of this. Starting from $b^{(t-1)}$ with a single $D \to B$ connection, the heuristic connects $A$ to $B$ upon observing $B$'s activation straight after activating $A$, and then $B$ to $C$ when $C$ activates shortly after.

2. **Time Sensitive (TS)** TS is like OO but with sensitivity to the expected cause–effect delays. It attributes activations to the (previous) event such that the respective delay would be most likely given the knowledge of the true causal delay distribution, and augments $b^{(t-1)}$ with an edge, if there is none yet, to form $b^{(t)}$. In the example (Figure 4b), $C$'s activation time is most consistent with $C$ being caused by the intervention on $A$, thus the model adds an $A \to C$ connection, rather than a $B \to C$ connection, going into $b^{(t+1)}$.

3. **Structure + Time Sensitive (STS)** STS is like TS, but it first checks if there is already an adequate explanation in the current model $b^{(t-1)}$. Concretely, it compares the likelihood of the most likely explanation that is already a cause in $b^{(t-1)}$ to the most likely explanation *overall* (i.e., the one selected by TS). Where these differ, it only adds an edge if the respective delay is substantially more likely than the delay implied by the best existing explanation in $b^{(t-1)}$, where we assume that "substantially more likely" means a likelihood ratio $> \frac{20}{1}$. Figure 4c gives an example. Unlike TS, this heuristic does not add an $A \to C$ connection going into $b^{(t+1)}$ because $C$'s activation can be explained well enough by the existing connection $D \to C$. While an $\mathbf{i}_A^{(2)} \to \mathbf{d}_B^{(1)}$ delay is slightly more probable than a $\mathbf{i}_D^{(1)} \to \mathbf{d}_B^{(1)}$ delay, the difference is not substantial enough to warrant the addition of another connection.

### Model comparison procedure

To compare the heuristics to participants' judgments, we simulated belief trajectories $\mathbf{b}$s for all the heuristics based on the evidence generated by all participants, starting each trial with an unconnected model at $t = 0$. For TS and STS, we assumed knowledge of true $\mu$, $\alpha$ and $w_S$ as participants had been trained on these during the instructions. We predicted participants' judgments based on what the simulated belief trajectories looked like at judgment time. We then assessed their accuracy in the task (e.g. the proportion of connections marked correctly) and accordance rate (the proportion of connections marked the same as the matched participant's). Additionally, we also compared participants to a *Random* baseline that marked a new random causal structure on every judgment, and an Ideal learner that always selects the max $P(M|\mathbf{d}_\tau; \mathbf{i}_\tau, \mathbf{w})$ according to the Bayesian inference model.

### Modeling results

The results of these simulations are reported in Table 1. Overall, *STS* was the most closely accordant with participants but individually participants were almost evenly split between STS and OO, both for all judgments and restricted to the final

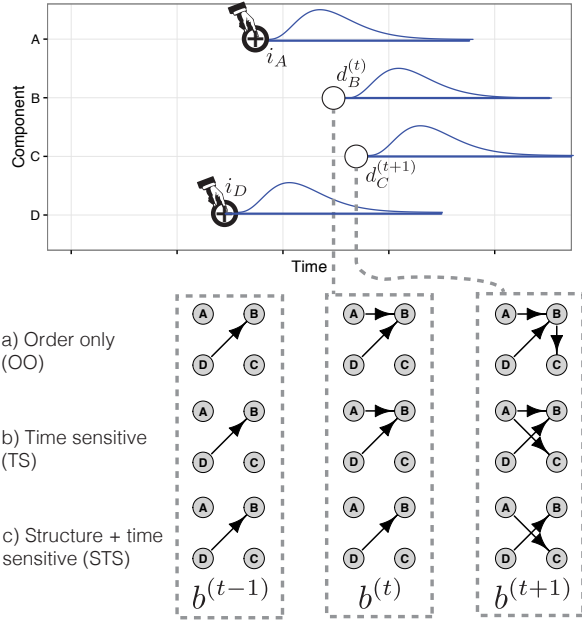Figure 4: Example where proposed heuristics' predictions diverge. $b^{(t-1)}$: the learner's belief at the start of the period depicted in the timeline. After observing $d_B^{(t)}$ the models update $b^{(t-1)}$ to form $b^{(t)}$. Then after observing $d_C^{(t+1)}$, they update to form $b^{(t+1)}$. Blue lines indicate probability density for cause–effect delays starting from each event, used to determine the most likely cause of each event (TS), and whether it is sufficiently more likely than any existing causes (STS).

judgments. Participants accuracy ($0.65 \pm 0.19$) was closest to that of the simplest heuristic OO. Mean participant accuracy by trial was correlated with that of all three heuristics $r_{OO} = .83, r_{TS} = 0.92, r_{STS} = 0.61$, but negatively correlated with Ideal judgments $r_{Ideal} = -.45$. Like participants but unlike the Ideal learner, all three heuristics were less accurate at cyclic than acyclic structures OO: $t(39) = 9.5, p < .001$, TS:$t(39) = 10.6, p < .001$, STS: $t(39) = 4.5, p < .001$.

## General Discussion

In our experiment, people used interventions to learn about the causal structure of devices whose dynamics unfolded in continuous time. As we predicted, cyclic structures were harder to learn than acyclic ones even though this was not reflected in the evidence available for an ideal learner, suggesting that the evidence produced by cyclic devices, involving many activations and potential causal paths, was harder for human learners to process. We found that the observed determinants of successful learning – equal spacing of interventions in time and a preference to intervene on root variables — made structure inference easier for a heuristic and

Table 1: Model comparison

| Model | Accuracy (%) | | Accordance (%) | | N best (/40) | |
| | All | Final | All | Final | All | Final |
|---|---|---|---|---|---|---|
| Random | 25.0 | 25.0 | 25.0 | 25.0 | 0 | 0 |
| OO | 66.2 | 64.7 | 67.2 | 64.9 | 16 | 17 |
| TS | 79.7 | 78.9 | 67.3 | 65.5 | 4 | 5 |
| STS | 87.9 | 90.9 | 69.3 | 69.2 | 15 | 13 |
| Ideal | 91.0 | 95.3 | 66.1 | 68.9 | 5 | 5 |

*Note:* "N Best" = the highest according model for each participant.

bounded learning system.

In light of this, we considered several heuristic learning models. Participants' judgments were best explained by assuming that they added connections to a single evolving candidate hypothesis as they observed events. Some subjects appeared to rely on a simple order heuristic (OO) whereas others displayed sensitivity to the delays between events (TS) and whether events were predicted by existing structure beliefs (STS). Participants rarely removed connections during the trials. Given more time to learn, however, it seems likely that they would also sometimes prune connections from their models — e.g., when events predicted by their current model repeatedly fail to occur. In general, positive testing of one's current hypothesis is an effective way for learners that are limited to a single global hypothesis to test its predictions against reality, and tune, refine, or or even abandon it, if necessary.

In sum, rather than grappling with an unmanageable space of possible structures and causal paths, participants seem to naturally follow Yogi Berra's advice: "You don't have to swing hard [to hit a home run]. If you got the timing, it'll go."

## References

Bonawitz, E. B., Denison, S., Gopnik, A., & Griffiths, T. L. (2014). Win-stay, lose-sample: A simple sequential algorithm for approximating Bayesian inference. *Cognitive Psychology*, *74*, 35–65.

Bramley, N. R., Dayan, P., Griffiths, T. L., & Lagnado, D. A. (2017). Formalizing Neurath's ship: Approximate algorithms for online causal learning. *Psychological Review*, *123*(3), 301–338.

Bramley, N. R., Gerstenberg, T., & Lagnado, D. A. (2014). The order of things: Inferring causal structure from temporal patterns. In *Proceedings of the 36th Annual Meeting of the Cognitive Science Society* (pp. 236–242).

Bramley, N. R., Gerstenberg, T., Mayrhofer, R., & Lagnado, D. A. (submitted). The role of time in causal learning.

Brehmer, B. (1992). Dynamic decision making: Human control of complex systems. *Acta Psychologica*, *81*(3), 211–241.

Buehner, M. J., & McGregor, S. (2006). Temporal delays can facilitate causal attribution: Towards a general timeframe bias in causal induction. *Thinking & Reasoning*, *12*(4), 353–378.

Coenen, A., Rehder, R., & Gureckis, T. M. (2015). Strategies to intervene on causal systems are adaptively selected. *Cognitive Psychology*, *79*, 102–133.

Fernbach, P. M., & Sloman, S. A. (2009). Causal learning with local computations. *Journal of Experimental Psychology: Learning, Memory & Cognition*, *35*(3), 678.

Gleick, J. (1997). *Chaos: Making a new science*. Open Road Media.

Greville, W. J., & Buehner, M. J. (2010). Temporal predictability facilitates causal learning. *Journal of Experimental Psychology: General*, *139*(4), 756–771.

Halpern, J. Y. (2016). *Actual causality*. MIT Press.

Lagnado, D. A., & Sloman, S. A. (2004). The advantage of timely intervention. *Journal of Experimental Psychology: Learning, Memory & Cognition*, *30*, 856–876.

Lauritzen, S. L., & Richardson, T. S. (2002). Chain graph models and their causal interpretations. *Journal of the Royal Statistical Society*, *64*(3), 321–348.

Nickerson, R. S. (1998). Confirmation bias: A ubiquitous phenomenon in many guises. *Review of General Psychology*, *2*(2), 175.

Pacer, M. D., & Griffiths, T. L. (2015). Upsetting the contingency table: Causal induction over sequences of point events. In *Proceedings of the 37th Annual Meeting of the Cognitive Science Society*.

Pearl, J. (2000). *Causality*. New York: Cambridge University Press (2nd edition).

Rehder, R. (2016). Reasoning with causal cycles. *Cognitive Science*, *to appear*.

Rottman, B. M., & Keil, F. C. (2012). Causal structure learning over time: observations and interventions. *Cognitive psychology*, *64*(1), 93–125.

Sloman, S. A., Love, B. C., & Ahn, W.-K. (1998). Feature centrality and conceptual coherence. *Cognitive Science*, *22*(2), 189–228.