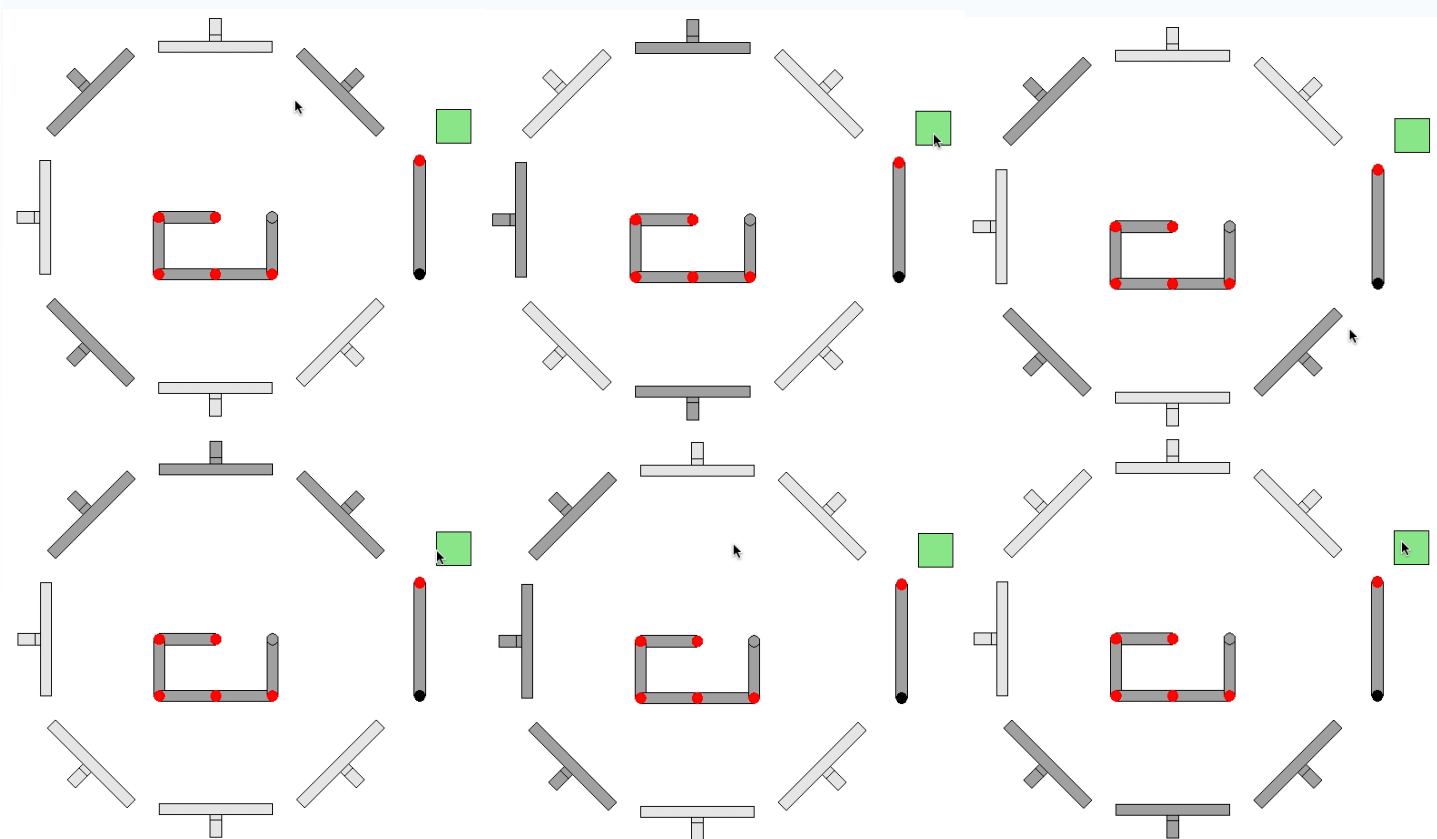


## INTRODUCTION

- Learning and storing abstract causal descriptions of the world enables generalization and transfer to new domains.
- Humans show a remarkable ability to form and utilize abstract causal structures to adapt and perform in novel domains.
- We task humans and agents to perform in a virtual “escape room” that requires reasoning about abstract causal structures and low-level properties of the scene. In this work, we model human causal learning as a two-component process:
  - A bottom-up associative account that links attributes of objects to causal effects.
  - A top-down causal structure account that encodes the latent structures most useful for the present task.
- The proposed model captures similar trends as human participants for multiple, suggesting human causal learning may rely on a synergy between a bottom-up associative learning scheme and a top-down structural learning scheme.

## METHOD & PROCEDURE

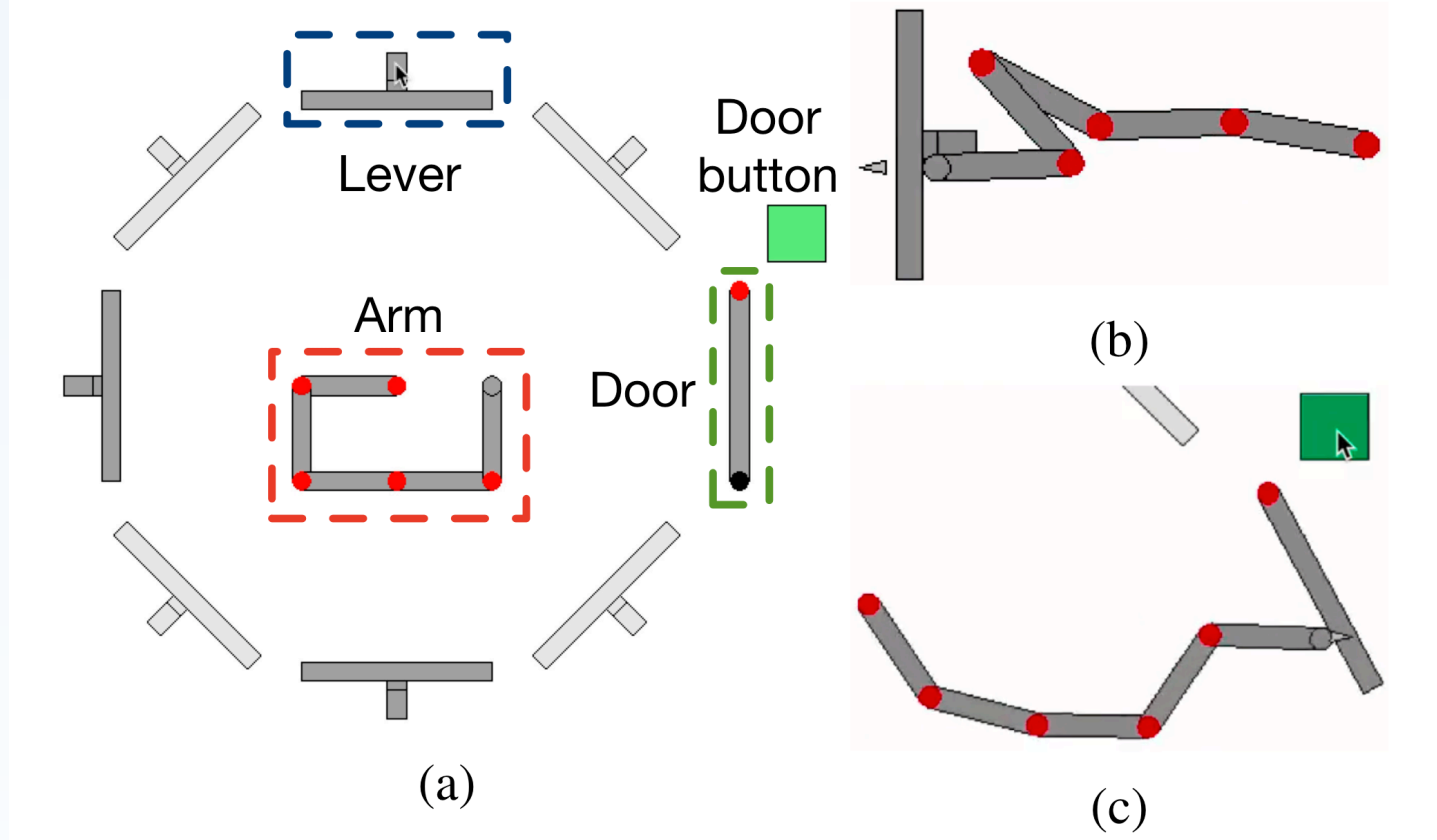
- We use the OpenLock task, initially presented in Edmonds et al., 2018. In the task, agents are required to “escape” from a virtual room by unlocking and opening a door.
- The door unlocks after manipulating the levers in a particular sequence (see Figure 1). Each room consists of seven levers surrounding a robotic arm that can *push* or *pull* on each lever.
- Agents observe the color of the levers and are expected to learn that grey levers—but not white levers—are always part of solutions in each room. Importantly, agents are tasked with finding all possible solutions for opening the door within a room.
- The mechanics underlying the environment obey one of two causal schemas: Common Cause (CC) and Common Effect (CE) (see Figure 3).



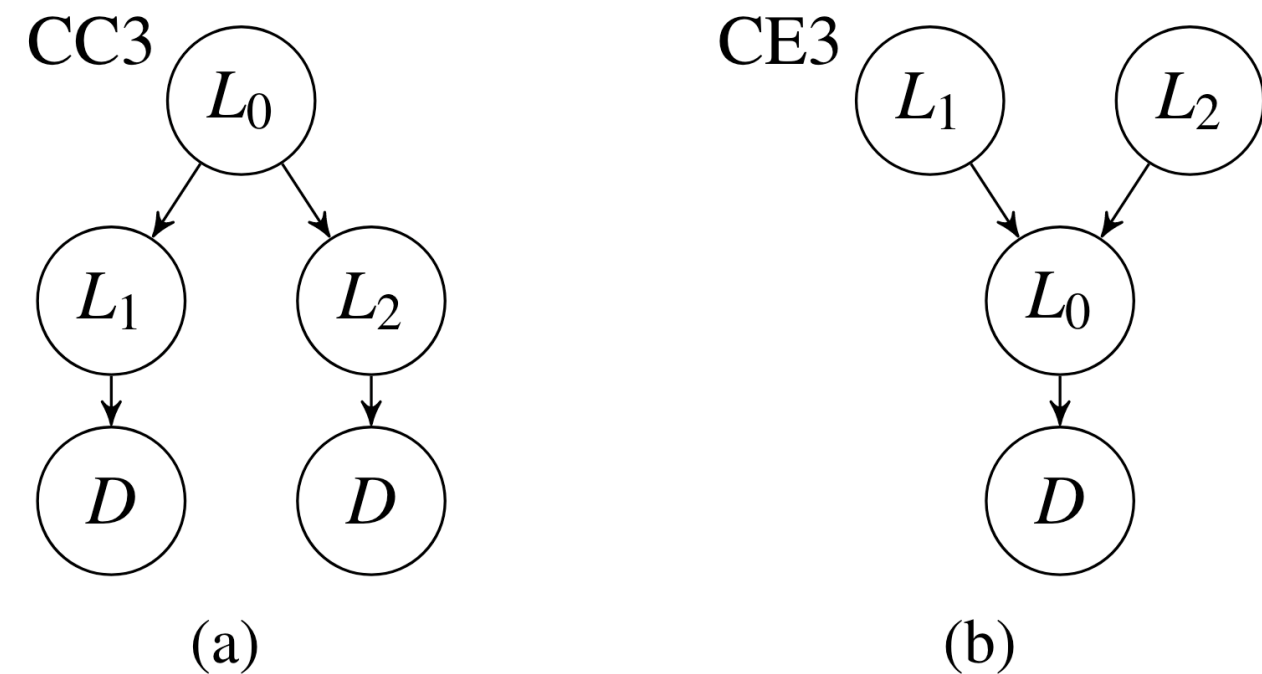
**Figure 1:** Six rooms used to ensure agents form abstract notion of task structure instead of overfitting to a specific room configuration.

## CAUSAL THEORY INDUCTION

- We approach the problem from the perspective of active causal theory learning, where we expect an agent endowed with no information to learn the underlying abstract mechanics and commonalities between environments through interaction.
- In this work, we adhere to two general principles of learning:
  - Causal relations induce state changes in the environment, and non-causal relations do not (referred to as our bottom-up  $\beta$  theory).
  - Causal structures that have previously been useful may be useful in the future (referred to as our top-down  $\gamma$  theory).



**Figure 2:** (a) Initial configuration of the room containing three active levers. The arm interacts with levers by pushing/pulling them outward/inward. Once the door is unlocked, the green button can be clicked to command the arm to push the door open. The black circle located opposite the door's red hinge represents the door lock indicator (present if locked, absent if unlocked). (b) Pushing on a lever. (c) Opening the door by clicking the green button.



**Figure 3:** Common Cause (CC) and Common Effect (CE) structures used in the OpenLock task, in which  $L_i$  indicates a lever in the scene, and  $D$  indicates the effect of opening the door.

**Attribute Learning:** Attributes provide time-invariant properties of an object; we learn which attributes are associated with causal events.

- For a particular causal chain, we want to assess the likelihood that the attributes ( $\phi_{ij}$ ) of the objects in the chain have been associated with causal events ( $\rho_i$ ) in the past:

$$p(\rho|c; \beta) = \prod_{c_i \in c} p(\rho_i|c_i; \beta) \quad p(\rho_i|c_i; \beta) \propto \prod_{\substack{\phi_{ij} \in s_i \\ s_i \in c_i}} p(\rho_i|\phi_{ij}; \beta)$$

**Schema Learning:** We utilize a Bayesian hierarchy, starting abstract structural schemas  $g^A$ , that encode abstract descriptions of the task.

- Using the Bayesian prior, we infer which instantiated schemas  $g^I$  are most likely to succeed, based on which abstract structures were useful in the past:

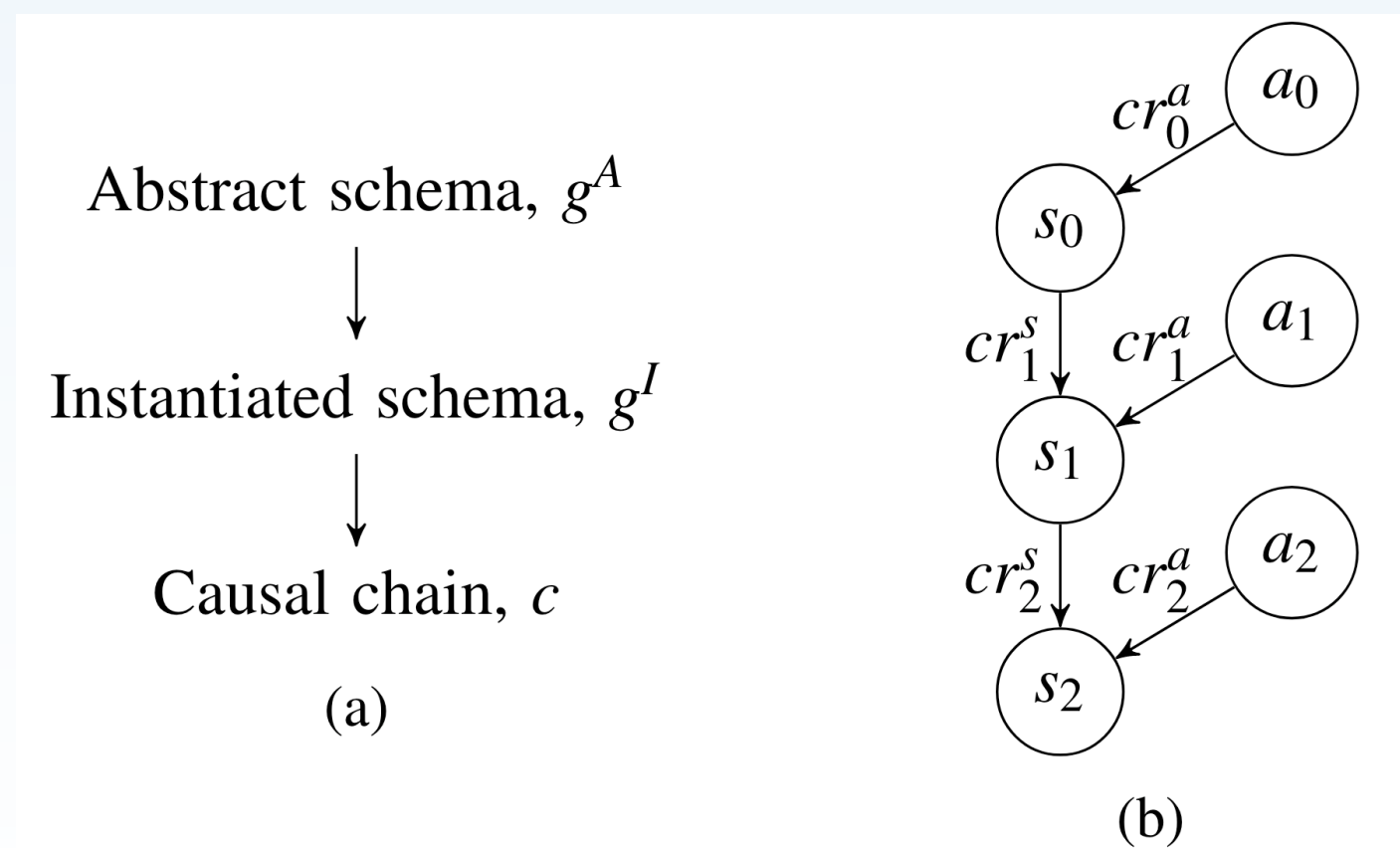
$$p(g^I|do(q); \gamma) = \sum_{g^A \in \Omega_{GA}} p(g^I|g^A, do(q))p(g^A; \gamma)$$

- Next we infer which chains are most useful based on which instantiated schemas are believed to be most useful:

$$p(c|do(q); \gamma) = \sum_{g^I \in \Omega_{GI}} p(c|g^I, do(q))p(g^I|do(q); \gamma)$$

- Finally, the agent makes decisions by combining the top-down schema reasoning (prior) and the bottom-up attribute learning (likelihood) to obtain a final posterior for a chain, and the agent executes the chain with the highest posterior at each time step:

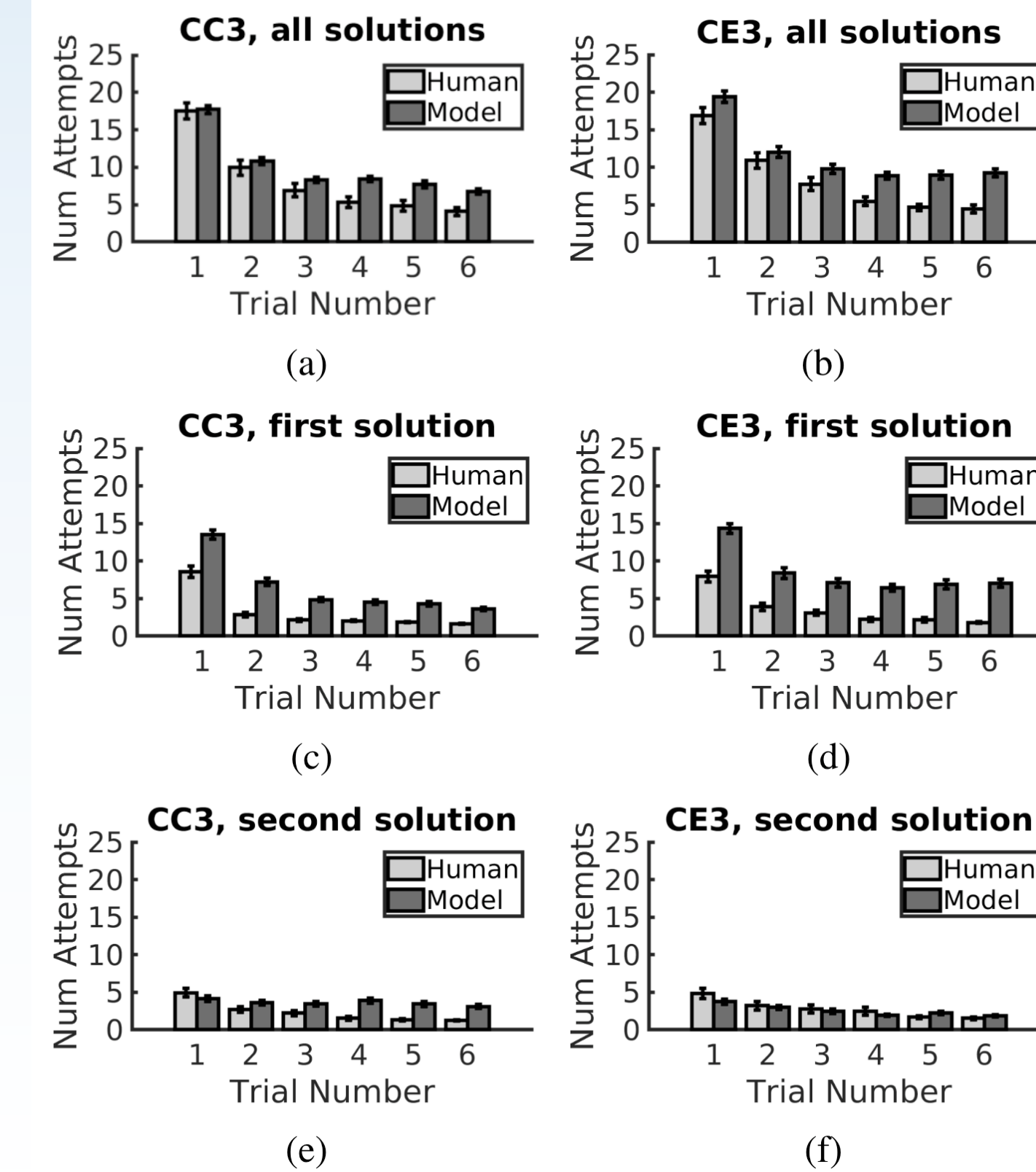
$$p(c|\rho, do(q); \gamma, \beta) \propto p(c|do(q); \gamma)p(\rho|c; \beta)$$



**Figure 4:** (a) An illustration of hierarchical structure of the model. A bottom-up associative learning theory,  $\beta$ , and a top-down causal theory,  $\gamma$ , serve as priors for the rest of the model. The model makes decisions at the causal chain resolution. (b) Atomic causal chain. The chain is composed by a set of subchains,  $c_i$ , where each  $c_i$  is defined by: (i)  $a_i$ , an action node that can be intervened upon by the agent, (ii)  $s_i$ , a state node capturing the time-invariant attributes and time-varying fluents of the object, (iii)  $cr_i^a$ , the causal relation between  $a_i$  and  $s_i$ , and (iv)  $cr_i^s$ , the causal relation between  $s_i$  and  $s_{i-1}$ .

## RESULTS

- The human results in Figure 5 demonstrate significant learning appeared to occur in the early trials for both the first and second solution.
- For participants who trained under a CC schema, attempts needed to find the *first* solution decreased significantly following both the first trial and second trial.
- For the *second* solution, the number of attempts needed decreased significantly following the first trial only.
- The model results show a similar trend as humans but with slightly worse performance in each trial.
- For the agent assigned to the CC condition, the number of attempts needed to find the first solution decreased significantly following the first trial and second trial.
- The CE agent required less attempts to find the first solution following the first trial only.



**Figure 5:** Comparison of human and model results for the common-cause CC3 condition and the common-effect CE3 condition. (a) and (b) compare the total number of attempts to find all solutions; (c) and (d) compare the number of attempts to find the *first* solution; (e) and (f) compare the number of attempts to find the *second* solution.

## CONCLUSION

- We showcase a hierarchical model based on associative learning and schema reasoning.
- Our model integrates two learning mechanisms:
  - A bottom-up theory that learns which attributes have causal associations in the environment
  - A top-down theory that learns useful abstract structures in the environment.
- Our agent chooses an intervention based on the posterior of causal chains and updates its model using the observed outcome of the intervention.
- Model results show that our hybrid agent is able to capture general trends observed in human participants and captures some of the statistical significance observed in human performance.
- These results suggest that human causal learning may consist of a mechanism that combines bottom-up associative learning with top-down reasoning about causal structure.

## DISCUSSION

**How can hypothesis space enumeration be avoided?**

- The spaces of  $\Omega_{g^A}$  and  $\Omega_{g^I}$  are enumerated in this work. Hypothesis space enumeration can quickly become intractable as problems increases. Future work will include examining how sampling-based approaches to iterative generate causal hypotheses.

**What are the other possibilities of bottom-up associative criteria?**

- Our method treats low-level attributes as the criteria for our bottom-up associative learning. However, other possibilities are equally valid; a modeler could pair attributes with specific actions and learn distributions of causal effects over this pairing.

## FUTURE DIRECTIONS

- The underlying computational framework presented here is broadly applicable outside of the OpenLock environment; it can be applied to any reinforcement learning environment where: (i) underlying dynamics are constrained by some causal structure; (ii) interactive elements have observable features which signal causal relevance; and (iii) physical locations of key elements change over time.
- We also hope to expand our model to account for more extreme observational changes. For example, what if levers could suddenly be rotated instead of pushed/pulled? What if new colors were introduced which provided further cues about causal relevance? And what if the environment began operating in a probabilistic fashion where levers may fail to actuate properly?