

Theory-based Causal Transfer: Integrating Instance-level Induction and Abstract-level Structure Learning

Mark Edmonds^{1,4}Xiaojuan Ma²Siyuan Qi^{1,4}Yixin Zhu^{2,4}Hongjing Lu^{2,3}Song-Chun Zhu^{1,2,4}UCLA Department of Computer Science¹UCLA Department of Statistics²UCLA Department of Psychology³International Center for AI and Robot Autonomy (CARA)⁴

INTRODUCTION

- Learning *transferable* knowledge across similar but different settings is a fundamental component to generalized intelligence.
- Our agent is endowed with two basic yet general theories for transfer learning:
 - A task shares a common abstract structure that is invariant across domains, and
 - The behavior of specific features of the environment remain constant across domains.
- RL agents showed poor ability transferring learned knowledge across different trials.
- The proposed model revealed similar performance trends as human learners, and more importantly, demonstrated transfer behavior across trials and learning situations.

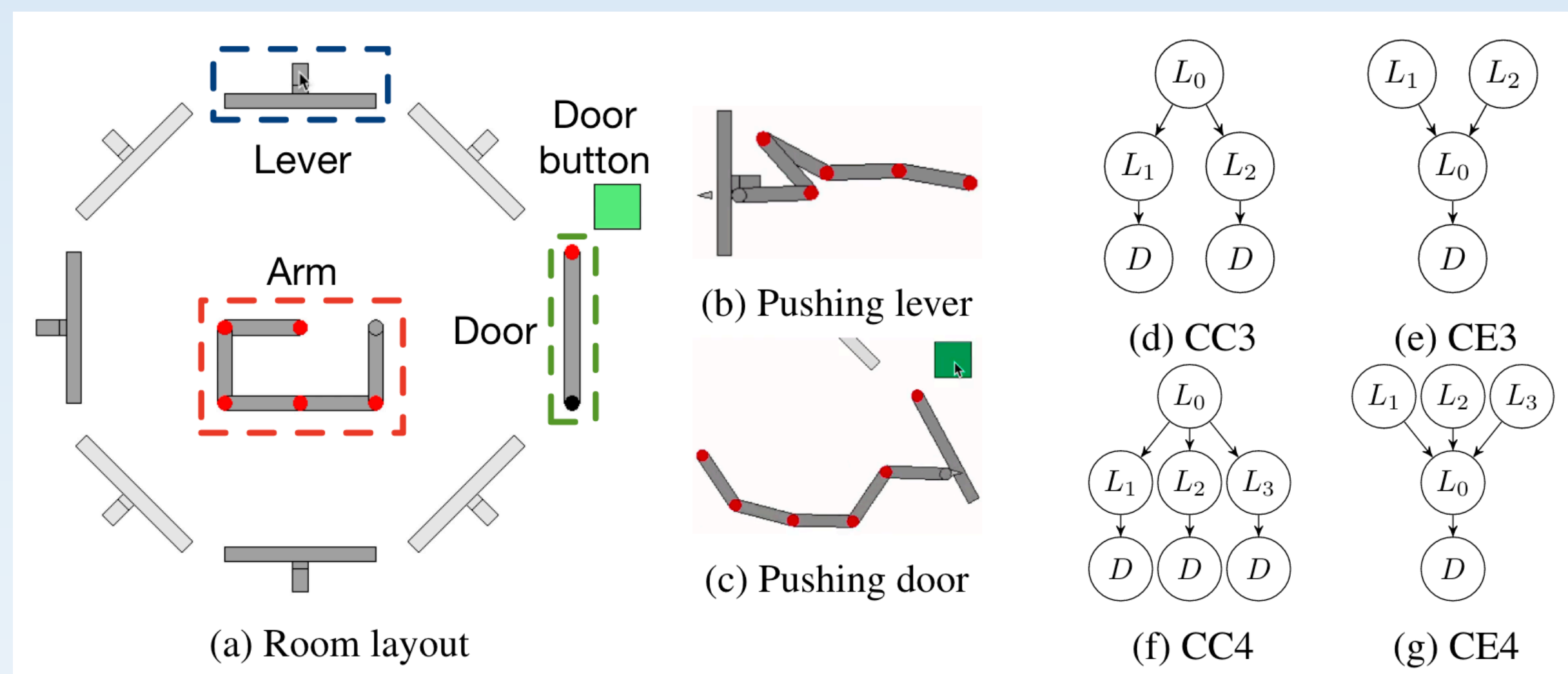


Figure 1: (a) Starting configuration of a 3-lever OpenLock room. The arm in the middle can interact with levers by either *pushing* outward or *pulling* inward. The door can be pushed only after being unlocked. The black circle on the door indicates whether or not the door is unlocked. (b) Pushing on a lever. (c) Opening the door. (d) CC3 causal structure. (e) CE3 causal structure. (f) CC4 causal structure. (g) CE4 causal structure.

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 β theory).
 - Causal structures that have previously been useful may be useful in the future (referred to as our top-down γ theory).

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

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

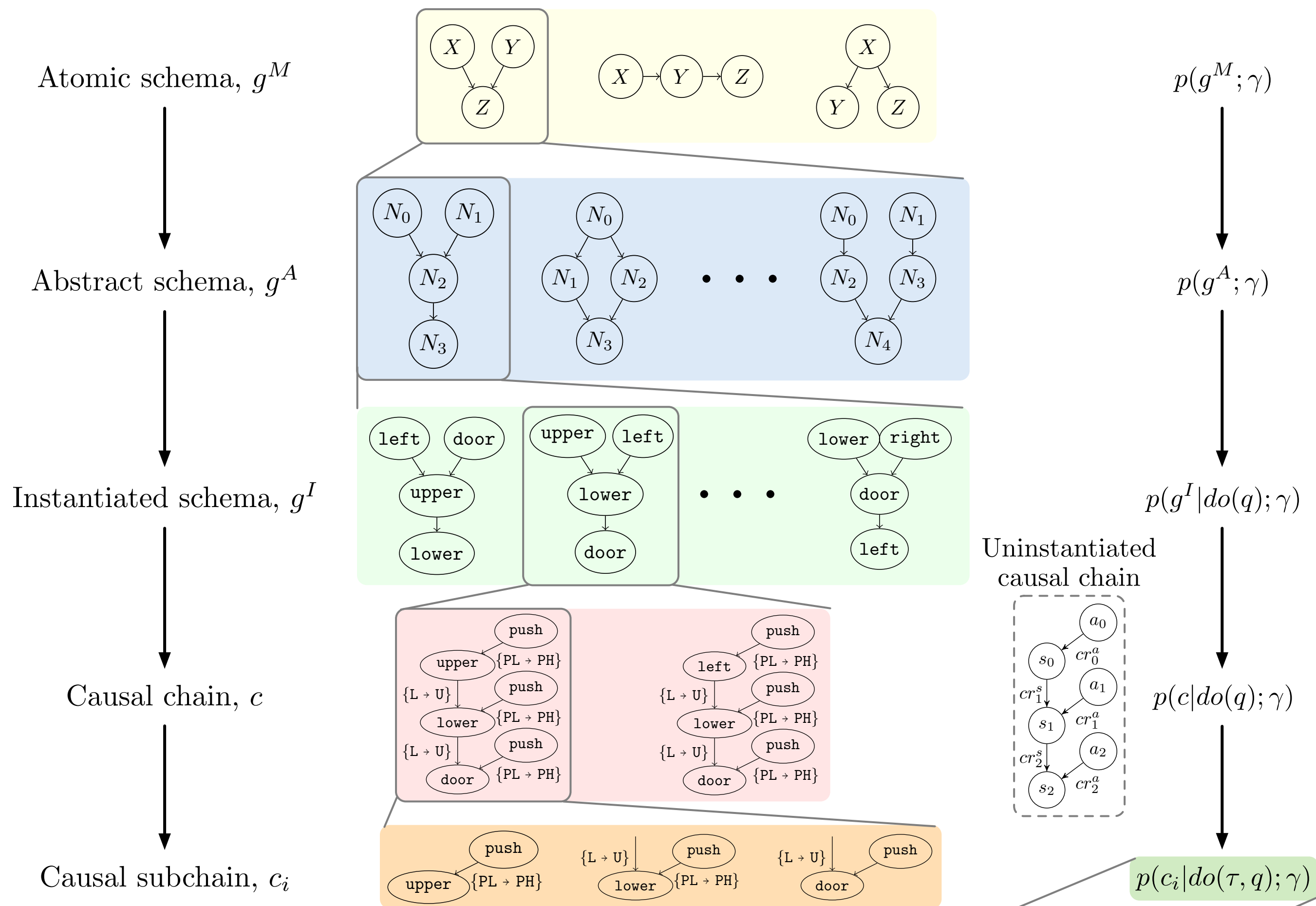
Intervention Selection: The agent marginalizes over the causal relations cr_i and states s_i to obtain a final, action-level term to select interventions:

$$p(a_i | \rho_i, do(\tau, q); \gamma, \beta) = \sum_{s_i \in \Omega_S} \sum_{cr_i^a \in \Omega_{CR}} \sum_{cr_i^s \in \Omega_{CR}} p(a_i, s_i, cr_i^a, cr_i^s | \rho_i, do(\tau, q); \gamma, \beta)$$

- The agent uses a model-based planner to produce action sequences capable of opening the door. The agent's final planning goal is

$$a_t^* = \arg \max_{a_i \in \Omega_{A^*}} p(a_i | \rho_i, do(\tau, q); \gamma, \beta)$$

(a) Abstract-level Structure Learning



(b) Subchain Posterior

$$p(c_i | \rho_i, do(\tau, q); \gamma, \beta) \propto p(\rho_i | c_i; \beta) p(c_i | do(\tau, q); \gamma)$$

(c) Instance-level Inductive Learning

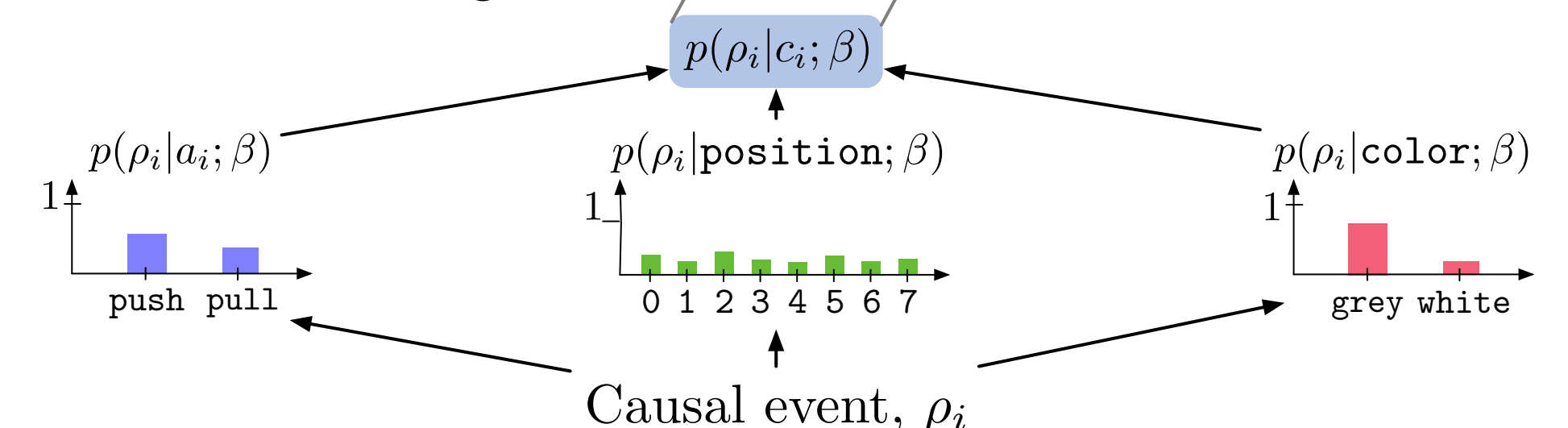


Figure 2: Illustration of top-down and bottom-up processes. (a) Abstract-level structure learning hierarchy. (b) The subchain posterior computed using the abstract-level structure learning and instance-level inductive learning. (c) Instance-level inductive learning.

RESULTS

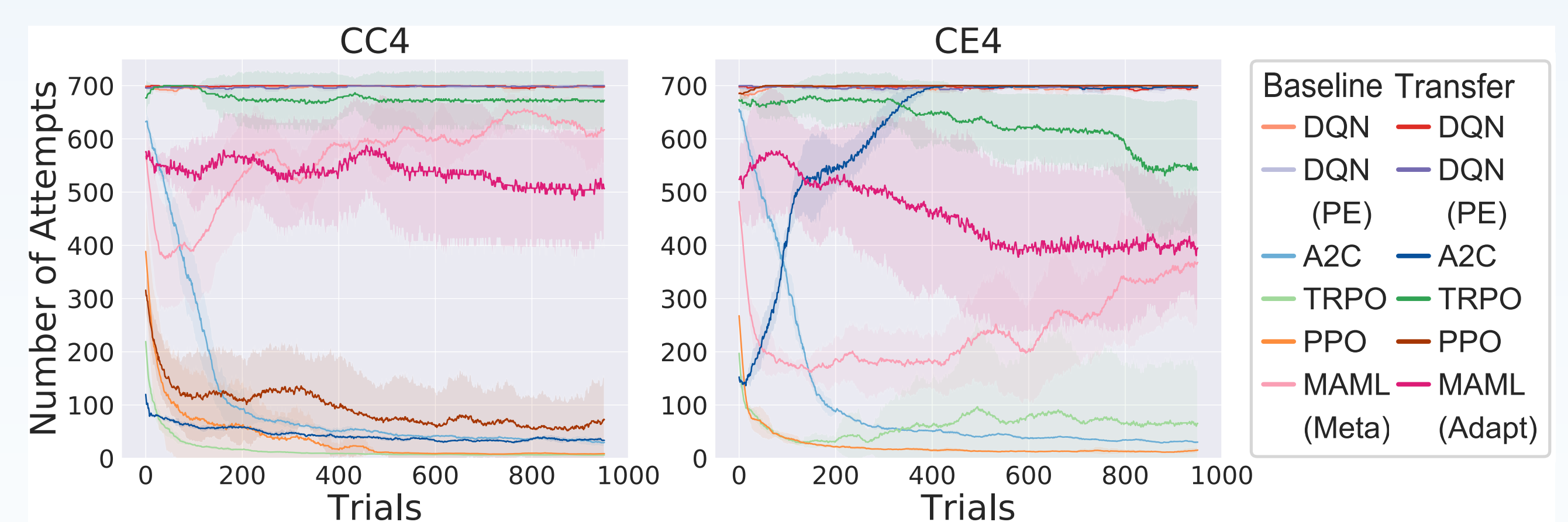


Figure 3: RL results for baseline and transfer conditions. Baseline (no transfer) results show the best-performing algorithms (PPO, TRPO) achieving approximately 10 and 25 attempts by the end of the baseline training for CC4 and CE4, respectively. A2C is the only algorithm to show positive transfer; A2C performed better with training for the CC4 condition. The last 50 iterations are not shown due to the use of a smoothing function.

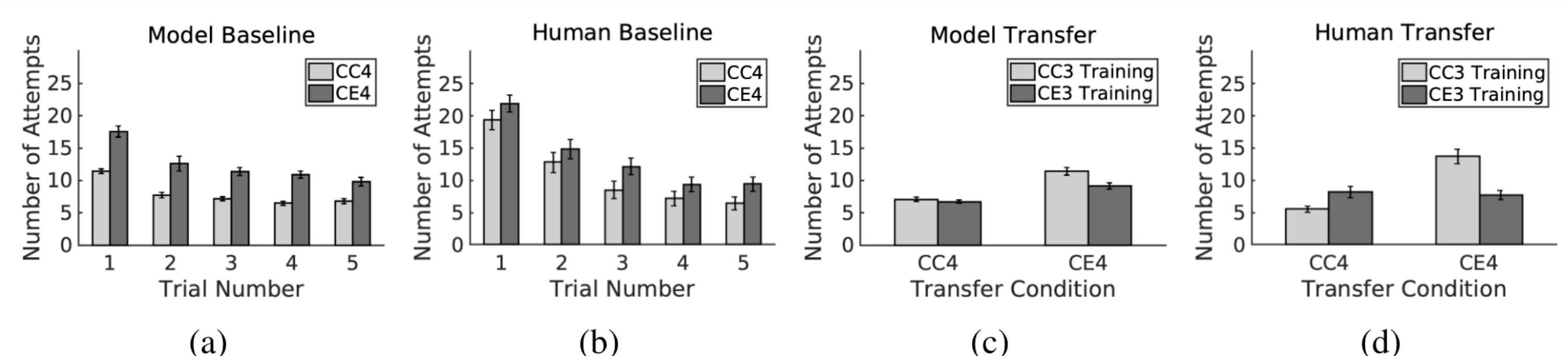


Figure 4: Results using the proposed theory-based causal transfer. (a) Proposed model baseline results for CC4/CE4. We see an asymmetry between the difficulty of CC and CE. (b) Human baseline performance from Edmonds et. al 2018. (c) Proposed model transfer results for training in CC3/CE3. (d) Human transfer performance from Edmonds et. al 2018.