# Using Causal Induction in Humans to Learn and Infer Causality from Video

Amy Fire (amy.fire@ucla.edu) and Song-Chun Zhu
Center for Vision, Cognition, Learning and Art, University of California, Los Angeles

**UCLA**

## INTRODUCTION

**Goal:** Computational model for the learning of causality from raw video

**Motivation:** Model inference processes

1. Answer why events occur
2. Correct misdetections and infer hidden/ambiguous objects/actions
3. Infer triggers, goals, and intents

## PERCEPTUAL CAUSALITY
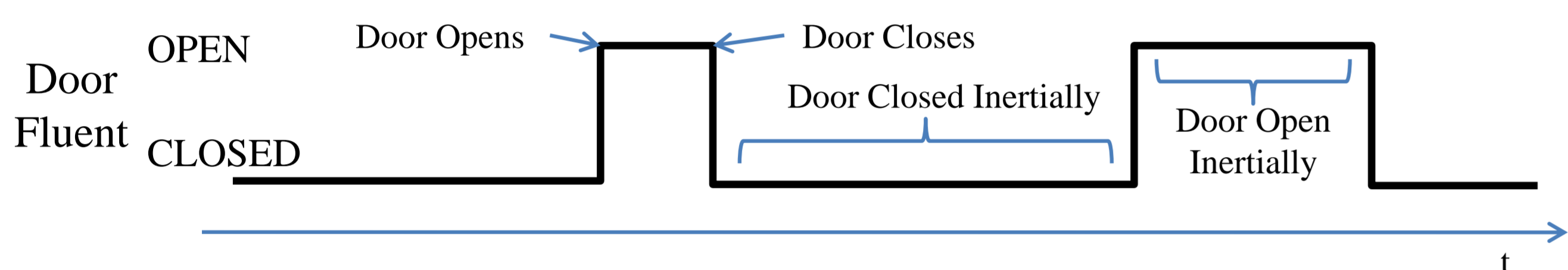
Infants use heuristics in judging causal relationships:

1. Agentive actions are causes
2. Measure co-occurrence between action $A_i$ and effect $\Delta F_j$

**cr** :

|  | $\Delta F_j$ Present | $\Delta F_j$ Absent |
|---|---|---|
| $A_i$ Present | $f_0$ | $f_1$ |
| $A_i$ Absent | $f_2$ | $f_3$ |

3. Temporal lag between the two is short
$$\text{Time(Action)} - \text{Time(Effect)} < \epsilon$$

4. Cause precedes effect
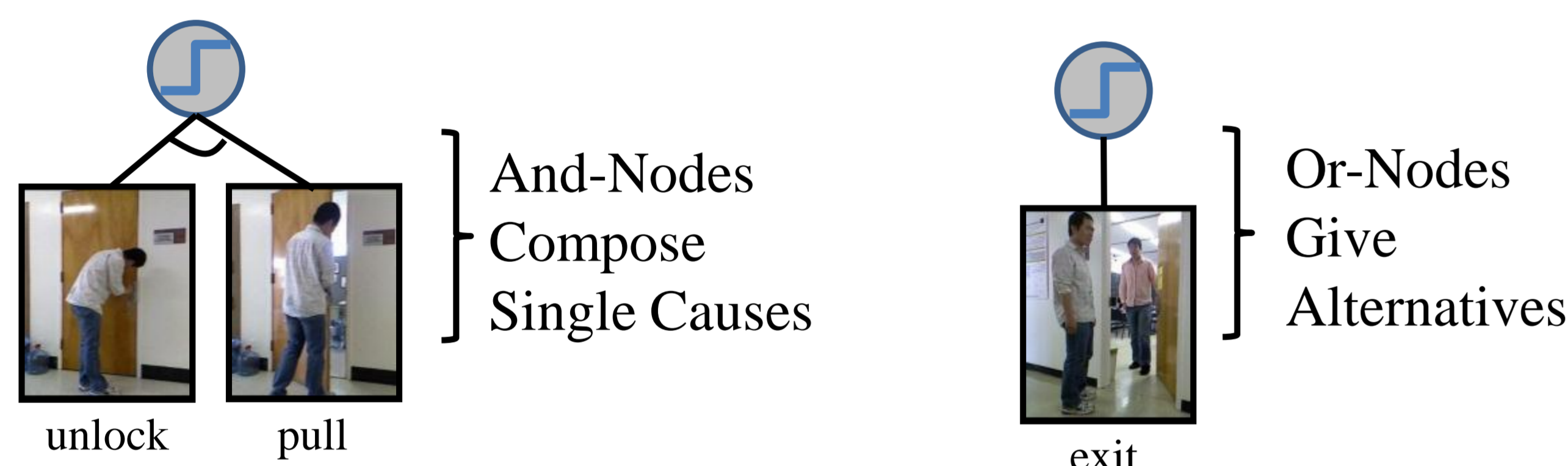$$\text{Time(Action)} - \text{Time(Effect)} > 0$$

To learn perceptual causality in video, we restrict co-occurrence of detected events and effects to these heuristics.
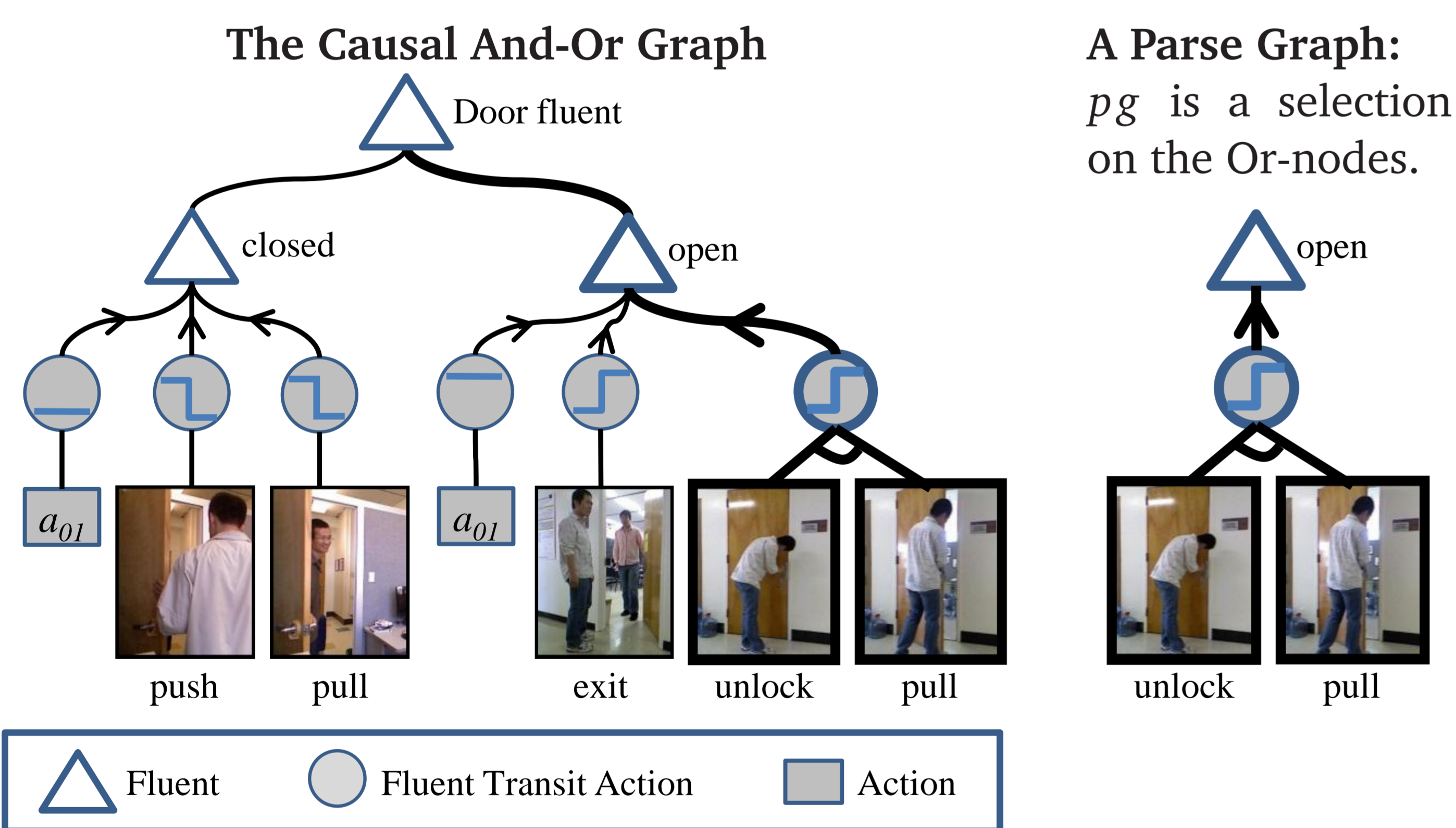
## A GRAMMAR MODEL FOR CAUSALITY

**Effects:** Fluents are time-varying statuses of objects.

Door Fluent: OPEN, CLOSED — Door Opens, Door Closes, Door Closed Inertially, Door Open Inertially, $t$

**Causes:** Actions suggest an And-Or representation.

And-Nodes Compose Single Causes — unlock, pull

Or-Nodes Give Alternatives — exit

**Pairing cause and effect:** Fluent changes are matched with corresponding causing actions. In the absence of change-inducing actions, fluent values are causally attributed to the *inertial action*, $a_{01}$.
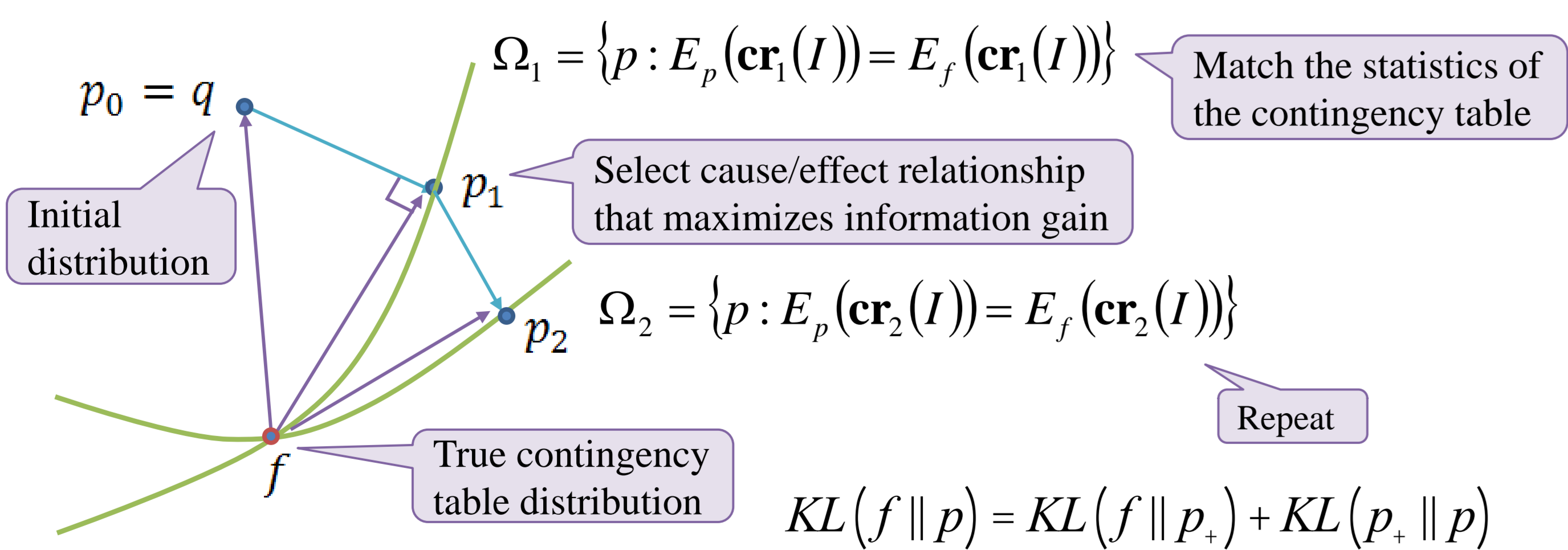
**The Causal And-Or Graph**

Door fluent — closed, open

$a_{01}$ — push, pull — $a_{01}$ — exit, unlock, pull

**A Parse Graph:** $pg$ is a selection on the Or-nodes.

open — unlock, pull

△ Fluent  ○ Fluent Transit Action  ▢ Action

**Probability on the C-AOG:** Given the video $I$,

$$\underbrace{P(pg_C|I)}_{\text{posterior}} = \underbrace{P(A_1,\ldots A_n|I)P(\Delta F_1,\ldots,\Delta F_m|I)}_{\text{likelihood}}\underbrace{\prod_{v \in V_C^{Or}}P(w(v))}_{\text{prior}}$$

- likelihood: the detection probabilities
- $V_C^{Or}$: the set of included Or-nodes in the causal explanation
- $w(v)$: the selected Or-branch
- prior: the switch probability on the Or-nodes

**Learning the C-AOG by model pursuit:** Incrementally pursue a model, adding a contingency table at each iteration by information projection.

$\Omega_1 = \{p : E_p(\mathbf{cr}_1(I)) = E_f(\mathbf{cr}_1(I))\}$ — Match the statistics of the contingency table

$p_0 = q$ — Initial distribution

Select cause/effect relationship that maximizes information gain

$\Omega_2 = \{p : E_p(\mathbf{cr}_2(I)) = E_f(\mathbf{cr}_2(I))\}$ — Repeat

True contingency table distribution

$$KL(f \parallel p) = KL(f \parallel p_+) + KL(p_+ \parallel p)$$

**Proposition:** Add the best action-fluent pair $(A_i, \Delta F_j)$:

$$\mathbf{cr}^* = \underset{cr}{\text{argmax}}(\text{Information Gain}) = \underset{cr}{\text{argmax}}(KL(\mathbf{f}\|\mathbf{h})),$$
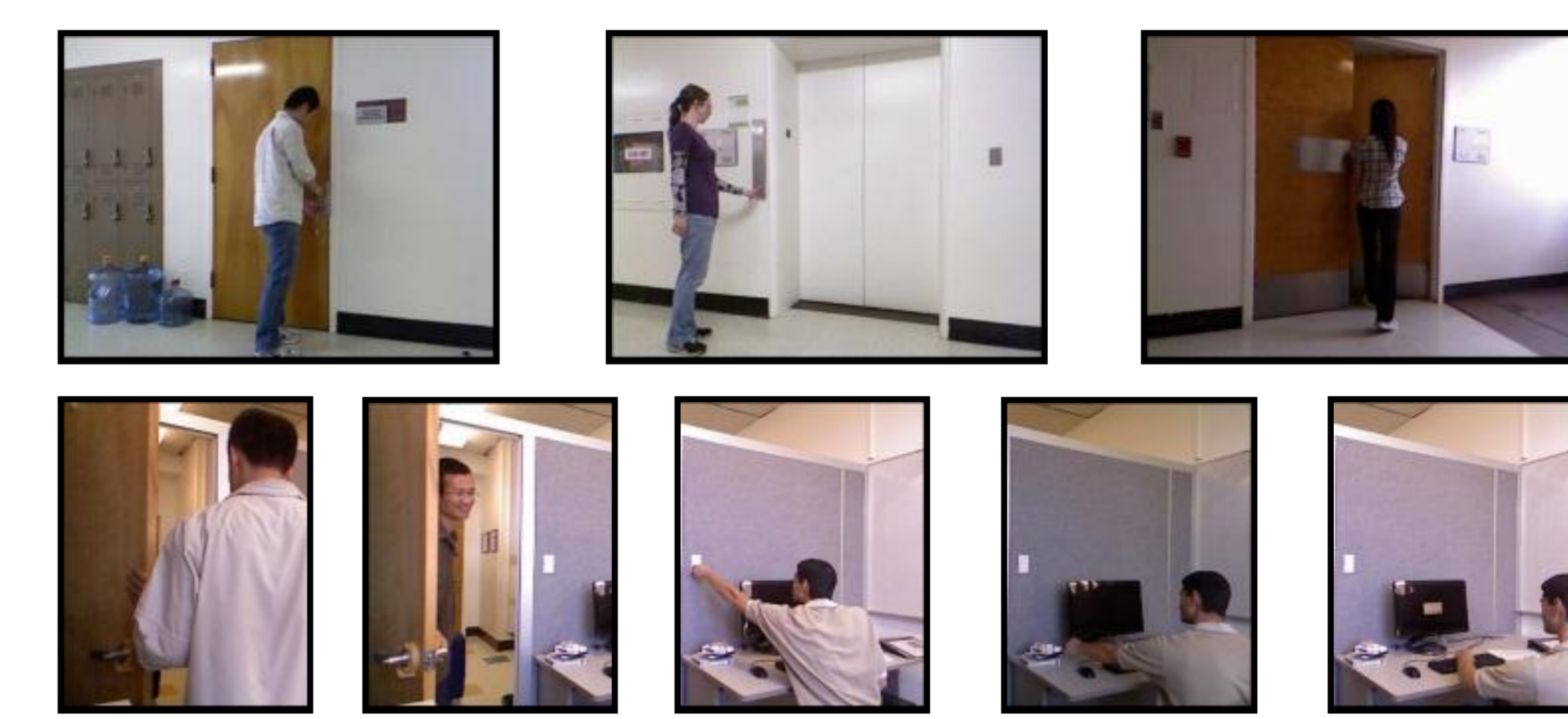
where $\mathbf{f}$ is the observed frequencies of $\mathbf{cr}$ and $\mathbf{h}$ is the expected contingency table predicted by the model $p$ in the current iteration
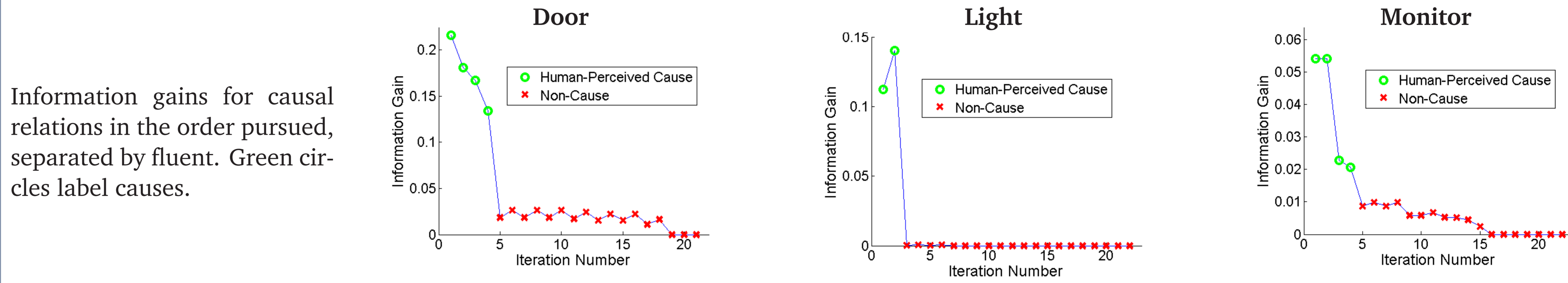
## EXPERIMENT 1: LEARNING CAUSALITY

**Goal:** Learn causal relationships between fluent changes and actions
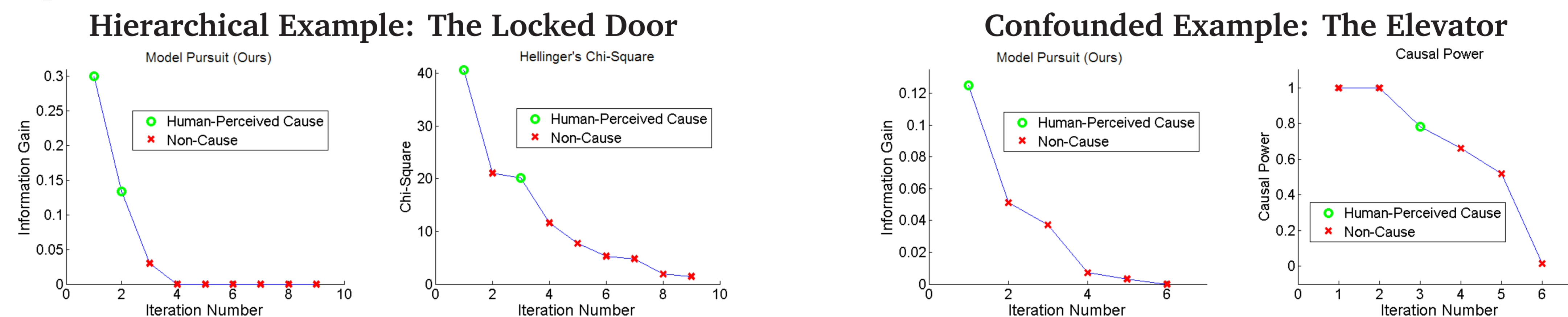
**Methods:**

- 120 minutes of video in office and hallway scenes
- 21 action categories, 8-20 instances of each
- Perfect action/fluent detection demonstrates learning
- Ground truth links known causing actions to their fluent effects
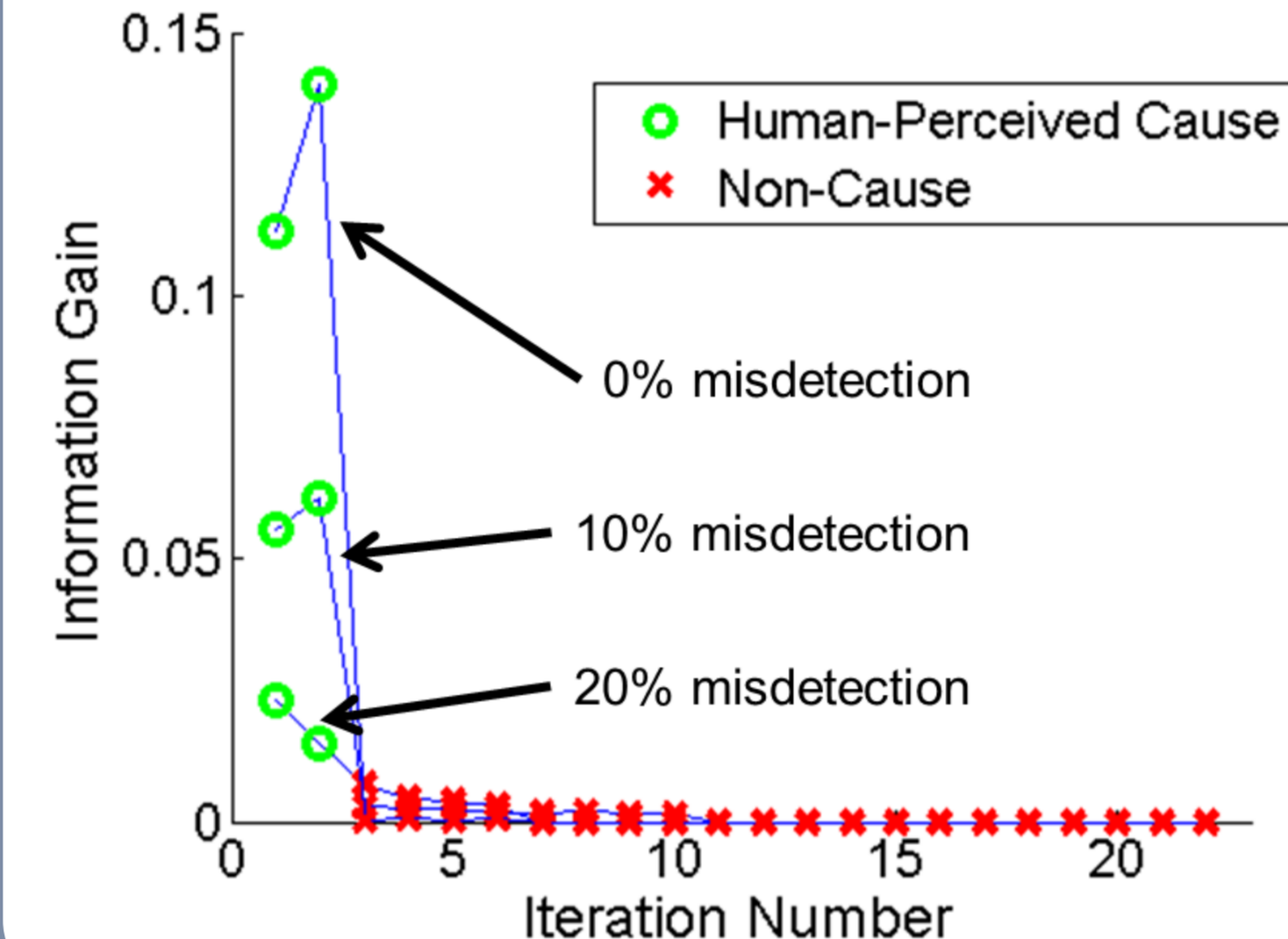
**Results: Correctly matching causal relations**

Information gains for causal relations in the order pursued, separated by fluent. Green circles label causes.


Door / Light / Monitor plots — Information Gain vs Iteration Number; Human-Perceived Cause (green circle), Non-Cause (red x)

**Comparisons:**

**Hierarchical Example: The Locked Door** — Model Pursuit (Ours), Hellinger's Chi-Square

**Confounded Example: The Elevator** — Model Pursuit (Ours), Causal Power

Our method acquires true causes before non-causes, outperforming Hellinger's Chi-Square and Causal Power.

**Noisy Data: Increasing Misdetections**


Information Gain vs Iteration Number; 0% misdetection, 10% misdetection, 20% misdetection

**Discussion:**

- Our method matches human perceptions in the presence of multiple confusing events.
- In the presence of confounders (the monitor), our method appropriately reduces clarity in the causal relationships.
- Our method incorporates dependencies in action hierarchies (the locked door).
- Our method places importance on quantity of hits (the elevator), accommodating the ambiguity important to vision.
- Clean detections are important to being able to learn causality.
- Limitation: Our methods are limited to pre-specified action and fluent categories so that appropriate detectors can be trained.

## EXPERIMENT 2: INFERENCE EXPERIMENT

**Goal:** Validate our model in the long-term reasoning task of inferring hidden fluent values

**Stimuli:**

- 20 minutes of hallway and office video
- 15 volunteer participants were shown the test video which paused at preset frames surrounding fluent changes or causing actions
- Fluents shown are either ambiguous or completely hidden


Frame Number (not to scale): 485, 800, 2500, 2575, 5535, 6915

**List of fluents**

Computer: ASLEEP/AWAKE
Monitor Display: ON/OFF
Monitor Power: ON/OFF
Cup: MORE/LESS/SAME
Water Stream: ON/OFF
Light: ON/OFF
Phone: ACTIVE/STANDBY
Trash Can: MORE/LESS/SAME
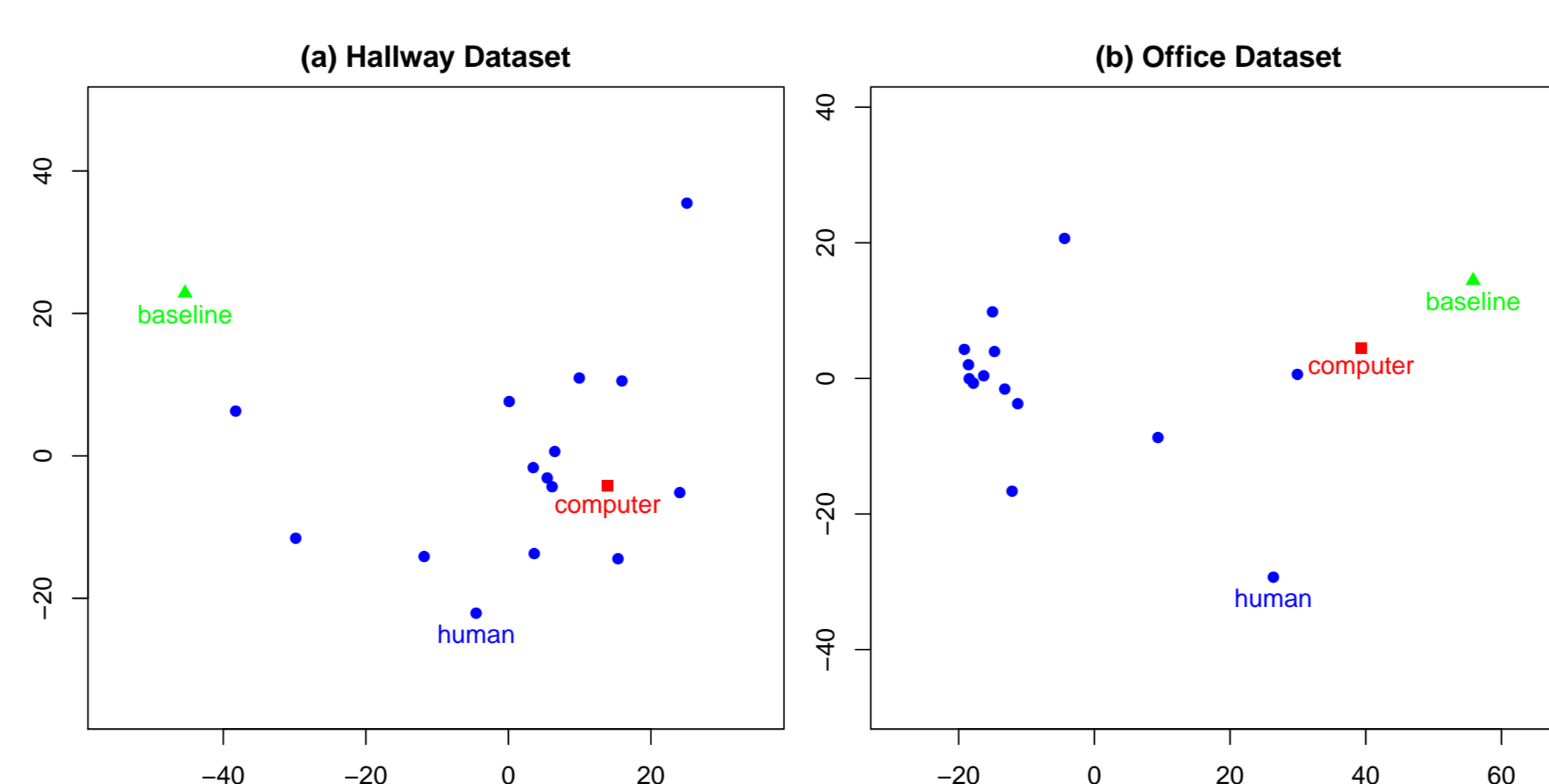Agent : THIRSTY/SATIATED
Agent: HAS_TRASH/NOT

**Reference estimates:**

- Baseline: 50/50
- Computer (our method): From video, actions are parsed using the Temporal And-Or Graph (right) and fluent changes are extracted using GentleBoost (below). These outputs are parsed with the Causal And-Or Graph.
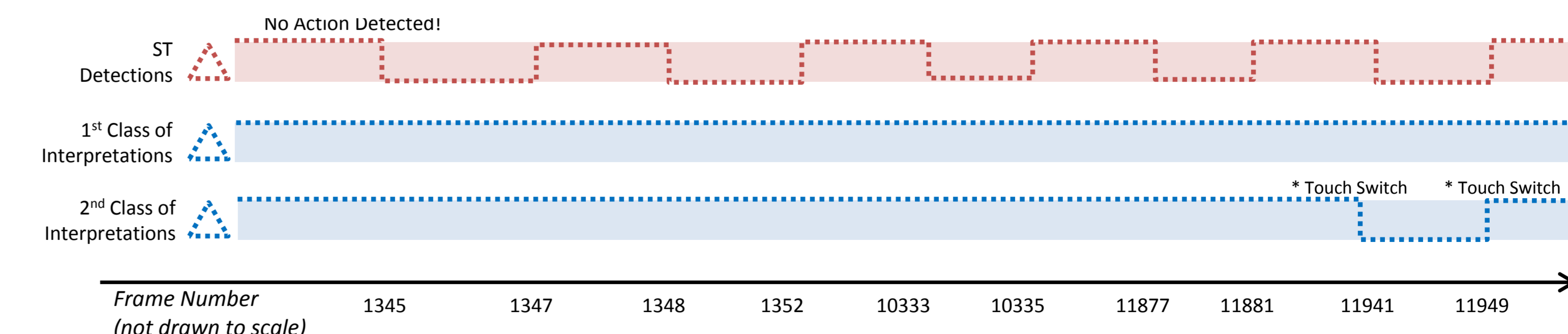
**T-AOG fragment**

parent(A) — is Part — $A$ — terminates — $R(A)$ — children(A)

Door open — Extract feature

door open, door closed, light on, light off, monitor off, monitor on — Non-maximum suppression

**Correcting Spatio-Temporal Detections:**

ST Detections — No Action Detected
1st Class of Interpretations
2nd Class of Interpretations — * Touch Switch, * Touch Switch
Frame Number (not drawn to scale): 1345, 1347, 1348, 1352, 10333, 10335, 11877, 11881, 11941, 11949

**Results:**

- MDS plots of fluent value estimates.


(a) Hallway Dataset, (b) Office Dataset — baseline, computer, human

**Discussion:**

- The Causal And-Or Graph smooths over misdetections in a way that is consistent with human responses
- The Causal And-Or Graph outperforms baseline
- Variation in human responses occurs due to different initializations and different variability thresholds