# Causal reasoning in a prediction task with hidden causes

Pedro A. Ortega, Daniel D. Lee, and Alan A. Stocker
University of Pennsylvania

# Motivation

- Humans guide decisions using **causal knowledge**.

- Causal knowledge predicts what the world does when we **interact** with it.

- Processing of causal information deeply embedded in animal cognition [1].

- Children develop causal understanding early on [2].

[1] Sloman, 2005; Blaisdell et al., 2006
[2] Gopnik et al. 2004; Meltzoff, 2007

# Motivation

Understanding how causal knowledge is

- **represented**,

- **learned**,

- and **used**

is currently **not well understood**.

[1] Sloman, 2005; Blaisdell et al., 2006
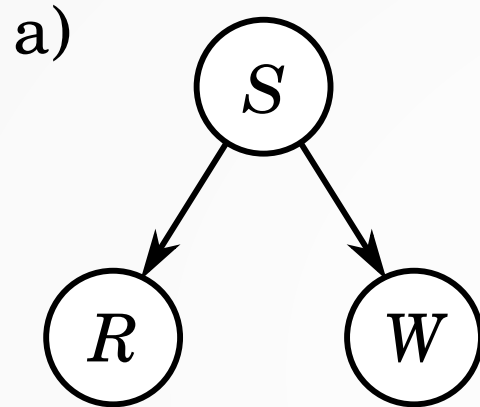[2] Gopnik et al. 2004; Meltzoff, 2007

# Causal theory of choice

- Humans infer **consequences** of their actions using **causal models** learned through experience [1].

- Causal knowledge is represented using **causal Bayes nets** [2].

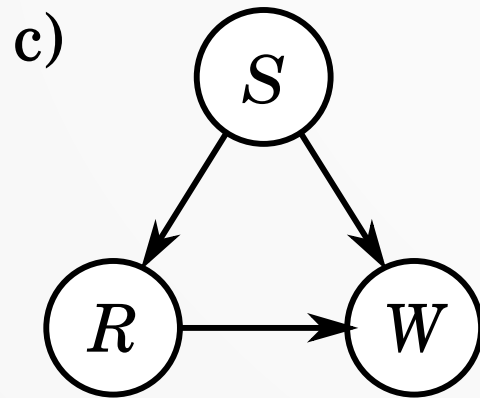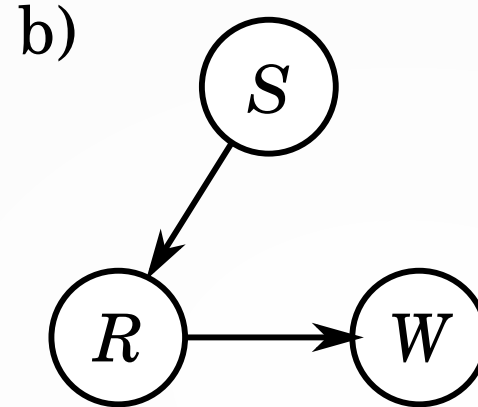[1] Hagmayer and Sloman, 2009
[2] Spirtes and Scheines, 2001; Pearl, 2009; Dawid, 2007
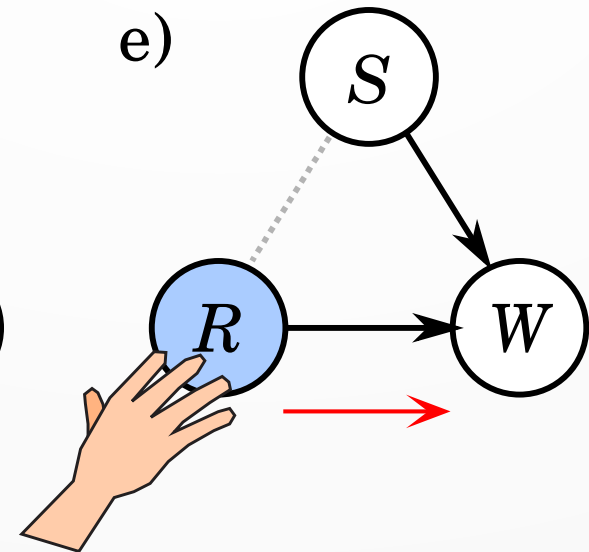
# Observations *vs.* Interventions



Common Cause Model

a)

Forward Model

b)

c)

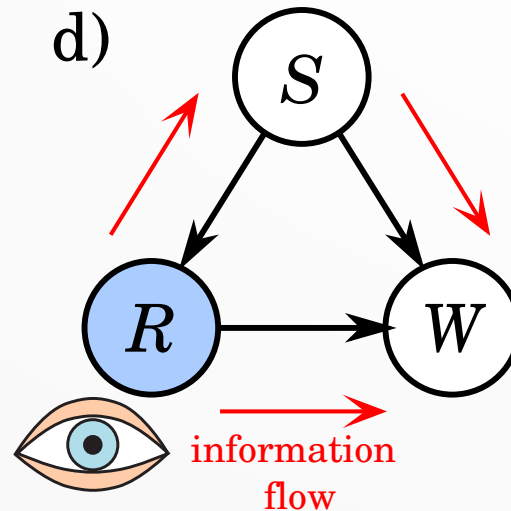Complex Model

d) information flow

e)

# Belief updates

- Observational:

$$P(W|R) = \sum_s P(W|S = s, R)P(S = s|R)$$

- Interventional:

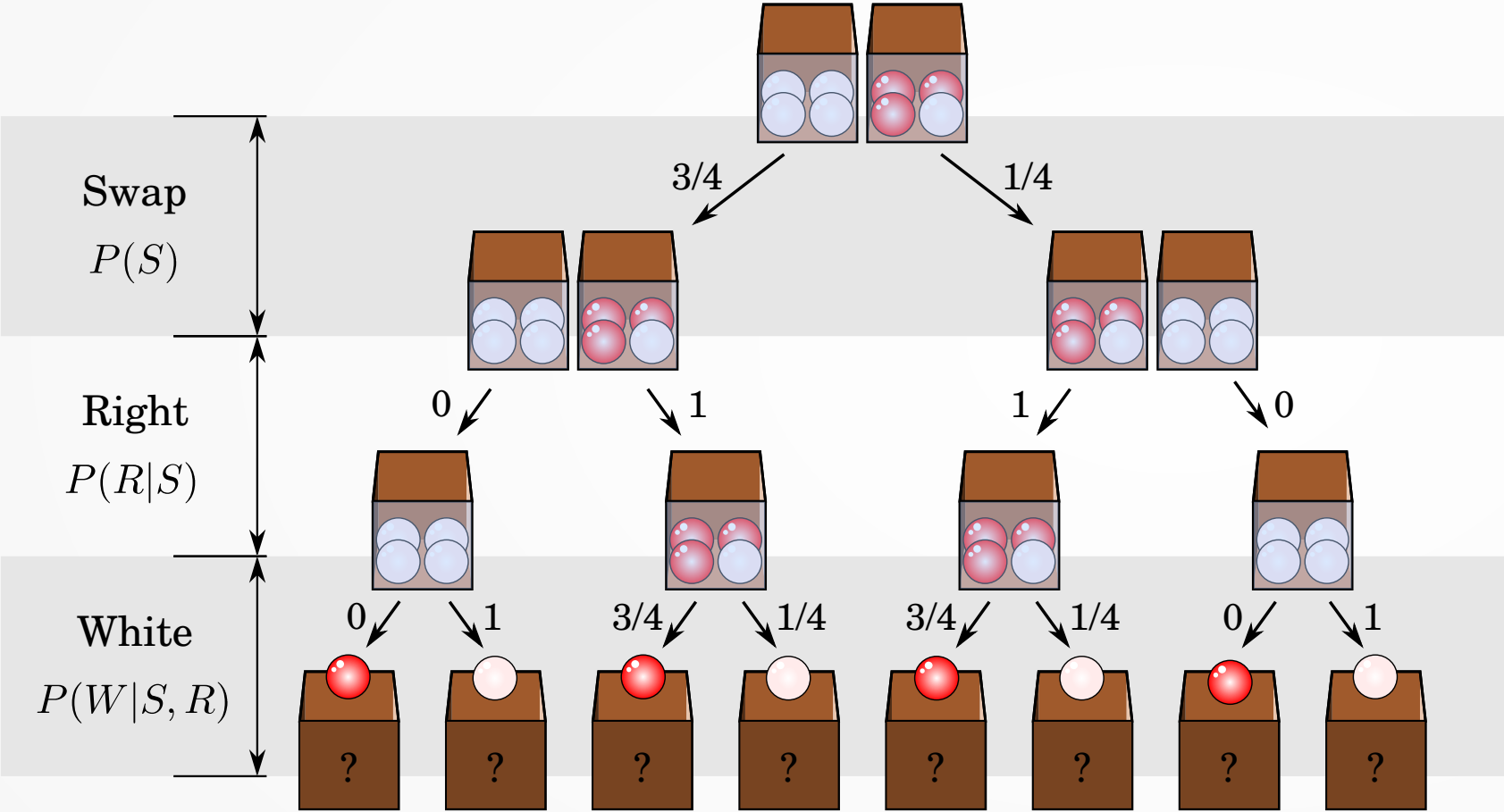$$P(W|\text{do}(R)) = \sum_s P(W|S = s, R)P(S = s)$$

# Questions

- Can humans learn and use **complex** causal structures?

- Hypothesis: Subjects **learn a complex causal dependency** (*i.e.* cause-effect relation) when they experience **both** the observational and interventional regimes.
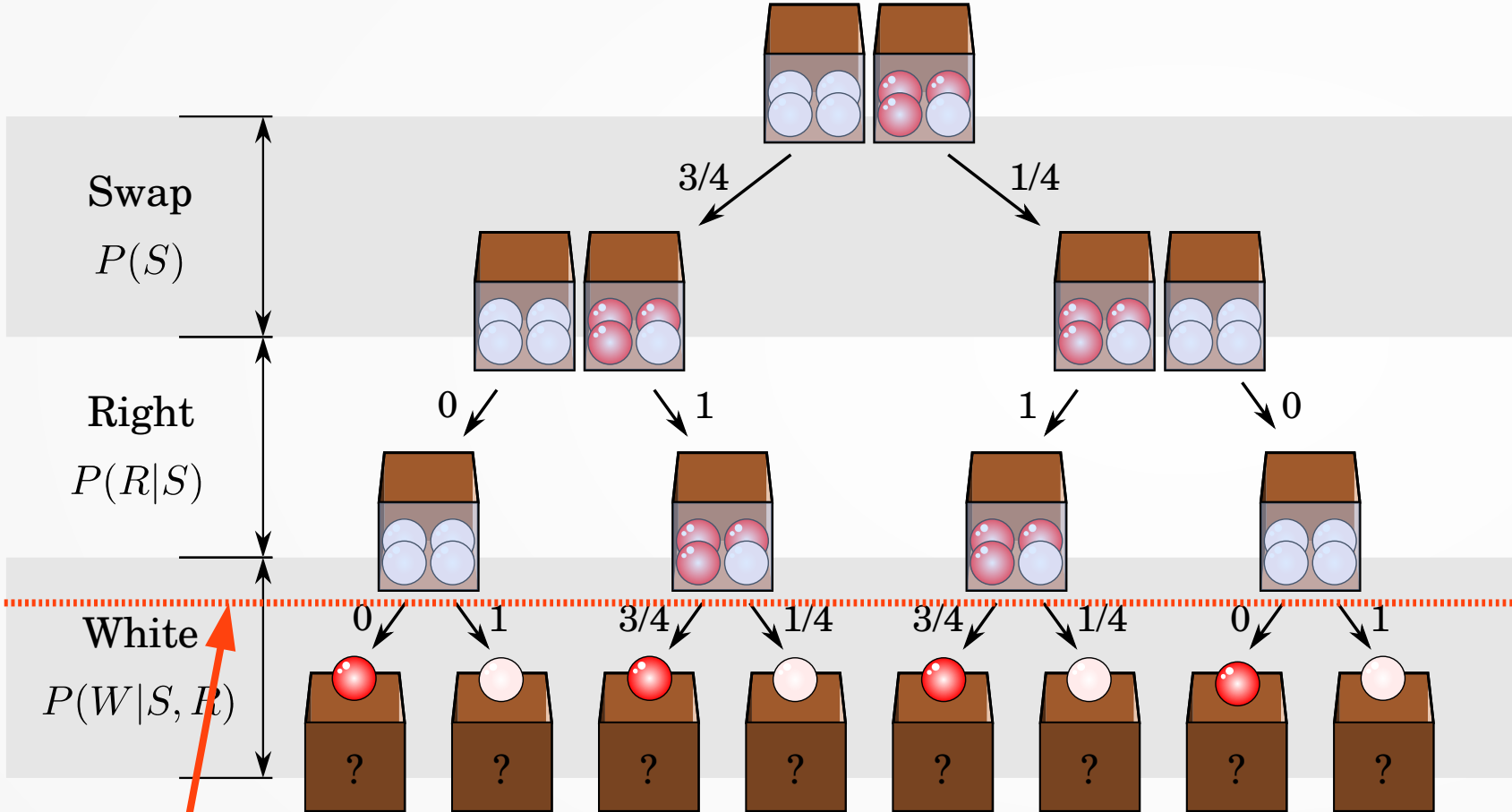
# Experimental method

- **Betting game** with hidden causes:

  - Two boxes with red and white balls.

  - Contents are **hidden**.

  - Bet on colour of **randomly drawn** ball.

- The causal structure is a complex model.

- Subjects play sequence of betting trials which they can **intervene** half of the time.

- We measure their **beliefs** and compare them to the model predictions.

# Betting game

# Betting game



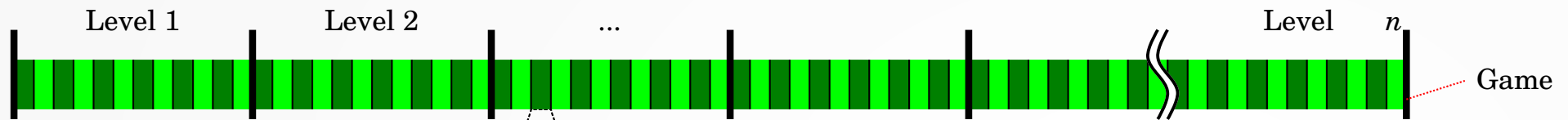Swap $P(S)$

Right $P(R|S)$

White $P(W|S,R)$

Place bet.

# Betting game



Swap $P(S)$

Right $P(R|S)$

White $P(W|S, R)$

3/4        1/4

0        1        1        0

0   1   3/4   1/4   3/4   1/4   0   1
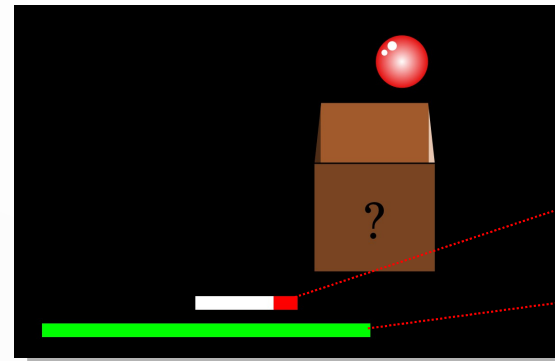
Place bet.

Choose left or right box.

11

# Game structure

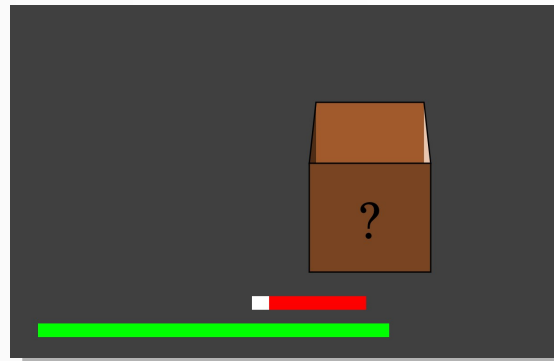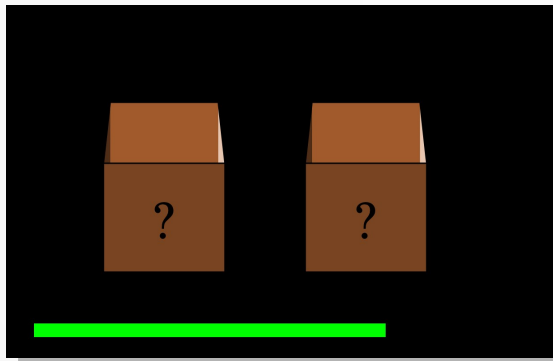- Subjects must complete 40 blocks (*levels*) of 10 trials each.

- They are allocated an **initial budget** at the beginning of each block.

- Each **bet reduces** the budget.

- Their goal is to **keep** as much as possible of the initial budget.

- If they reach zero, they **must repeat** the block.

# Game structure
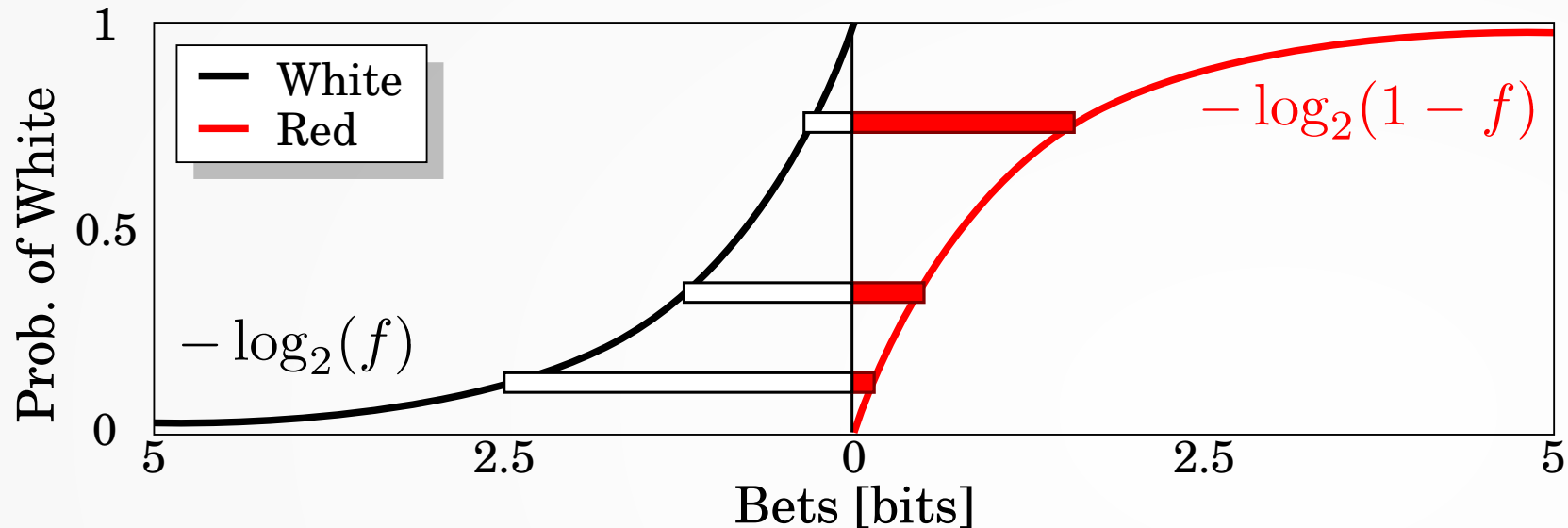


Level 1  Level 2  ...  Level *n*

Game

Trials

betting bars

remaining budget

a)  b)  c)

# Betting mechanism



- *Log-loss scoring rule* encourages reporting **true beliefs** [1].

- Allows measuring beliefs on a **trial-by-trial** basis.

- Confident bets are **too risky**.

- Initial budget **prevents conservative** guesses.

[1] Dawid, 2006; Bickel, 2007

# Training & test games

| Game | Levels | Transparent | Intervention |
| --- | --- | --- | --- |
| Training 1 | 10 | yes | no |
| Training 2 | 10 | yes | yes (50%) |
| Test | 40 | no | yes (50%) |

- We trained subjects on two simplified games:
  - Training 1 familiarises subjects with betting scheme.
  - Training 2 teaches the causal structure.
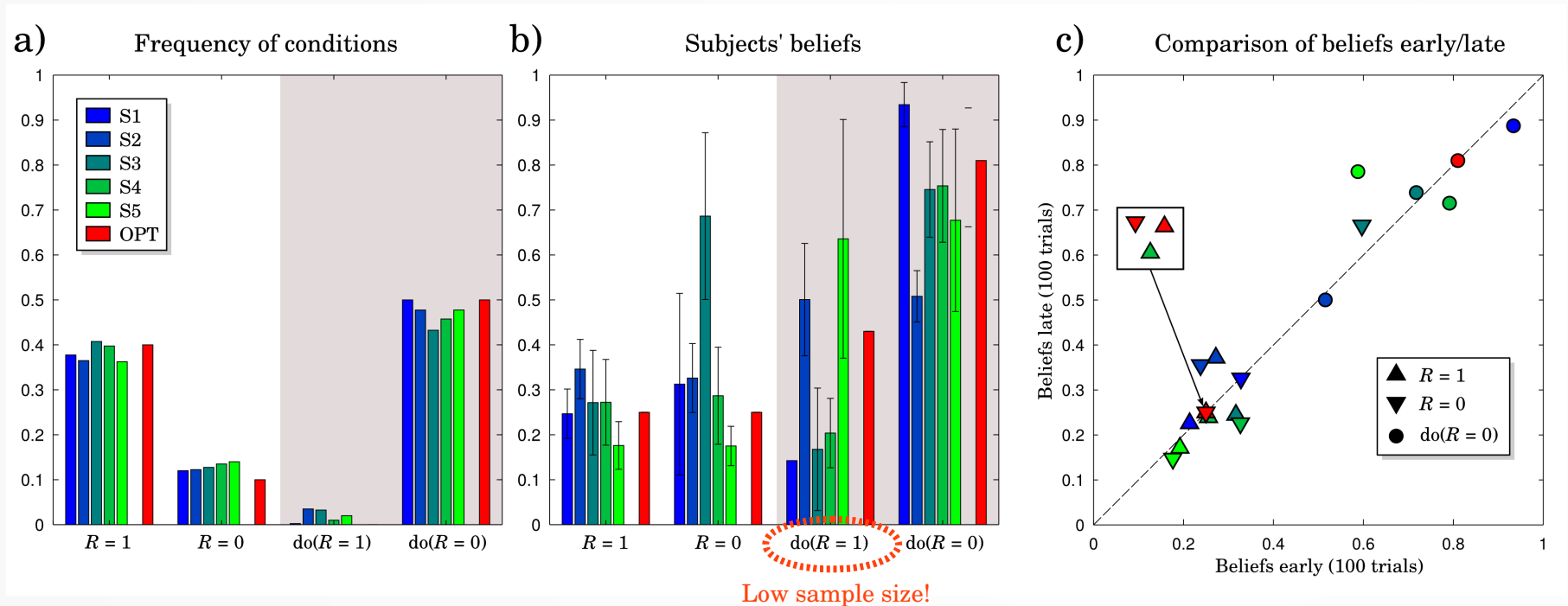
# Summary of experimental method

- **Betting** optimally requires:

  - **learning** the trial parameters (statistics and causal structure),

  - **marginalising** over then hidden causes,

  - and **distinguishing** between actions and observations.

- **To train** the subjects:

  - we let them play two short **training games**,

  - where the **contents** of the boxes were **visible** at all times,

  - and where we let them **experience each condition** half of the time.

- **To test** whether they use causal reasoning:

  - we measure their **predictive beliefs** about the ball's colour,

  - and **compare** them to the **model** predictions.
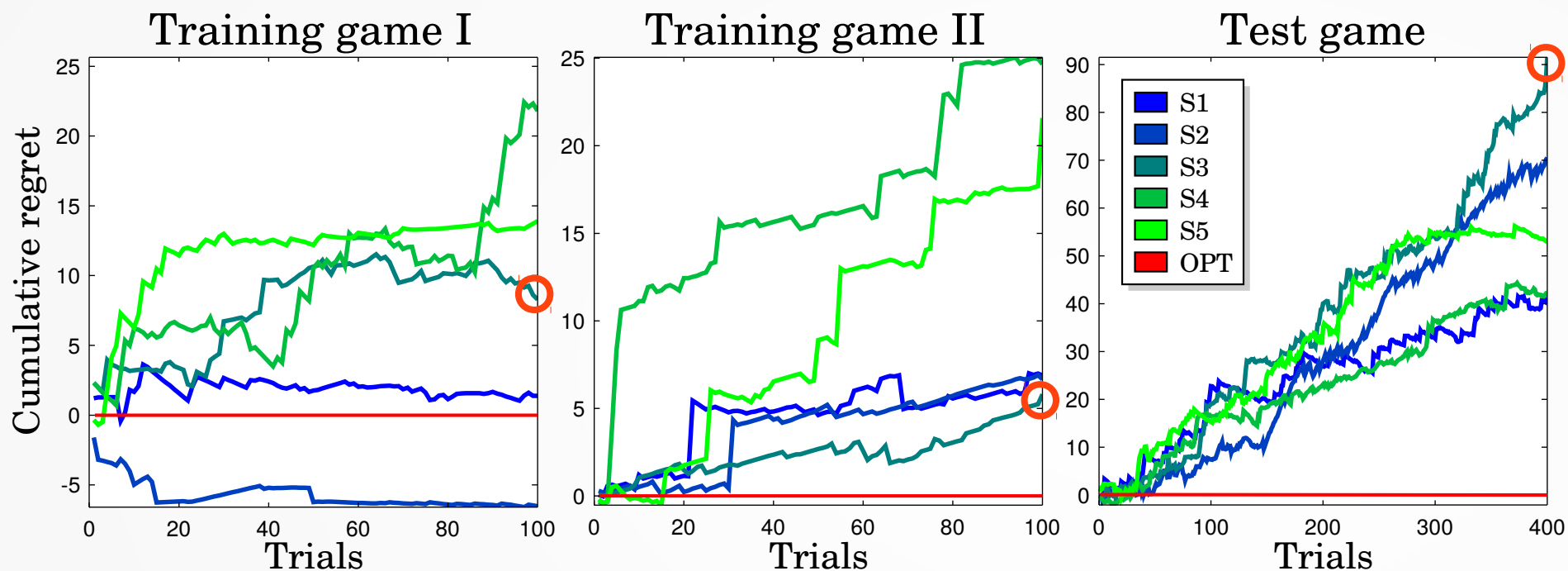
# Data collection

- Subjects: Five (UPenn) students (S1-S5).

- The training and test games were played in a **single session** (< 90 mins), totalling more than **600 trials**.

- Were not told statistics nor causal structure.

- Were told that all trials had identical statistics & causal structure; and the differences between games.

- $10 for participation + $10 for completion.

# Final prediction probabilities



a) Frequency of conditions

b) Subjects' beliefs

Low sample size!

c) Comparison of beliefs early/late

- 4 out of 5 learned to predict correctly **right from the start**.

- Combines expected utility, Bayes, and causality.

- S3 treated every condition as interventional.

# Learning curves



- Cumulative regret = performance - optimal.
- Smaller slope = better; negative curvature = learning.
- Training games: learning is very quick (< 40 trials).
- Test game: little to no learning—but positive slope: noisy beliefs?
- Curiosity: **S3** performs pretty well during the training games: smaller hypothesis space?

# Summary of results

- Excepting S3, all the subjects made bets that were **consistent** with the **causal model**'s predictions.

- Hence, they **induced** the causal model, **marginalised** over hidden causes, and **distinguished** between actions and observations.

- Crucially:

  - **absence** of learning during test game,

  - and **uselessness** of regime distinction during training games,

  suggest that subjects could **spontaneously** supply "regime indicators" to their experience.

# Conclusions

- Subjects can **learn complex** causal structures— it appears to be **sufficient** to let them experience both regimes.

- Subjects can use **causal deductive** reasoning.

- Subjects appear to **spontaneously** tag experience as either interventional or observational, even though they **do not need** to so to perform well.