

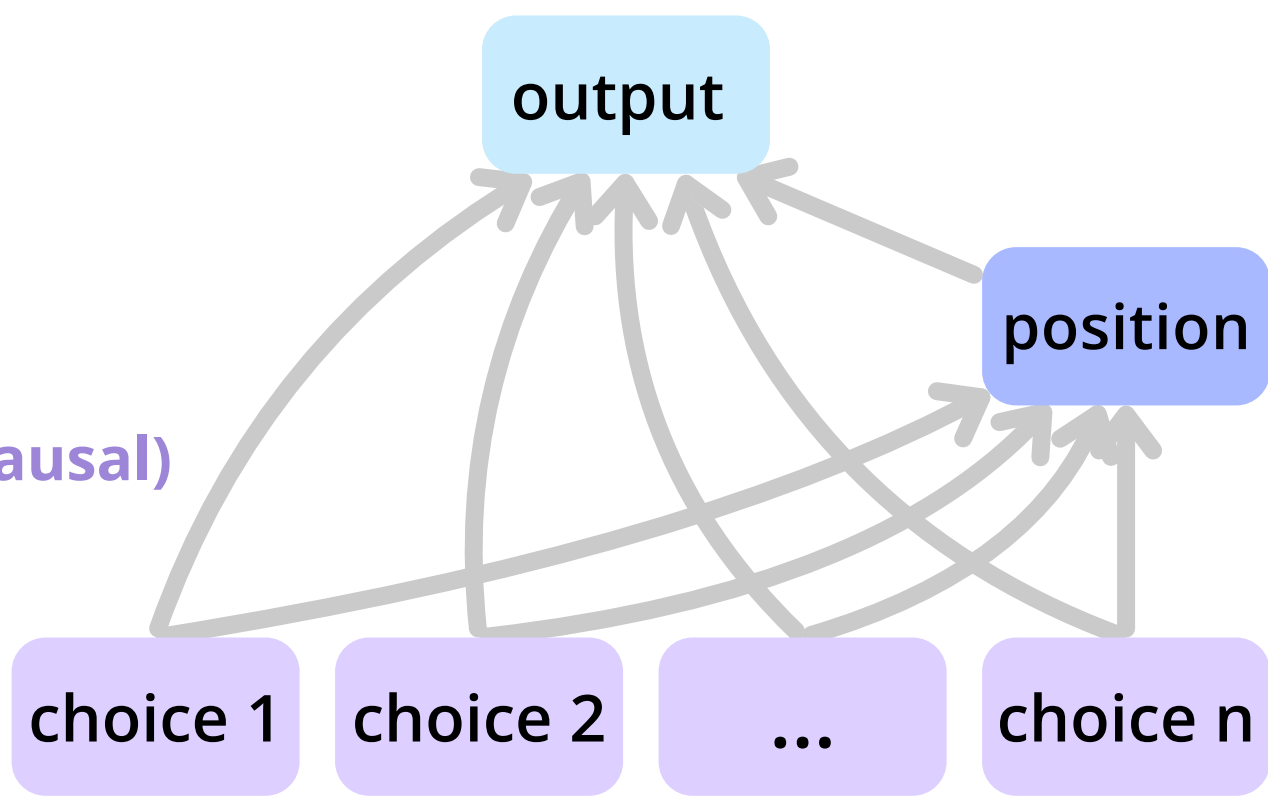
Interp Finding: Causal Mechanisms of MCQA

Output variable

Causal variables for OOD prediction

Causal variables

Background (non-causal) variables



Task: Predict OOD Behaviors on MMLU

Find the degree for the given field extension $Q(\sqrt{2}, \sqrt{3}, \sqrt{18})$ over Q .

ID Scenario	OOD Scenario
A. 0	Alpha. 0
B. 4	Bravo. 4
C. 2	Charlie. 2
D. 6	Delta. 6

Answer: B. **Answer: *Delta.***

Methods: Abstraction → Prediction

The model solves a task successfully → it likely implements a **systematic solution**, i.e. a **causal mechanism**

Abstract the high-level causal model from ID examples that model correctly solves

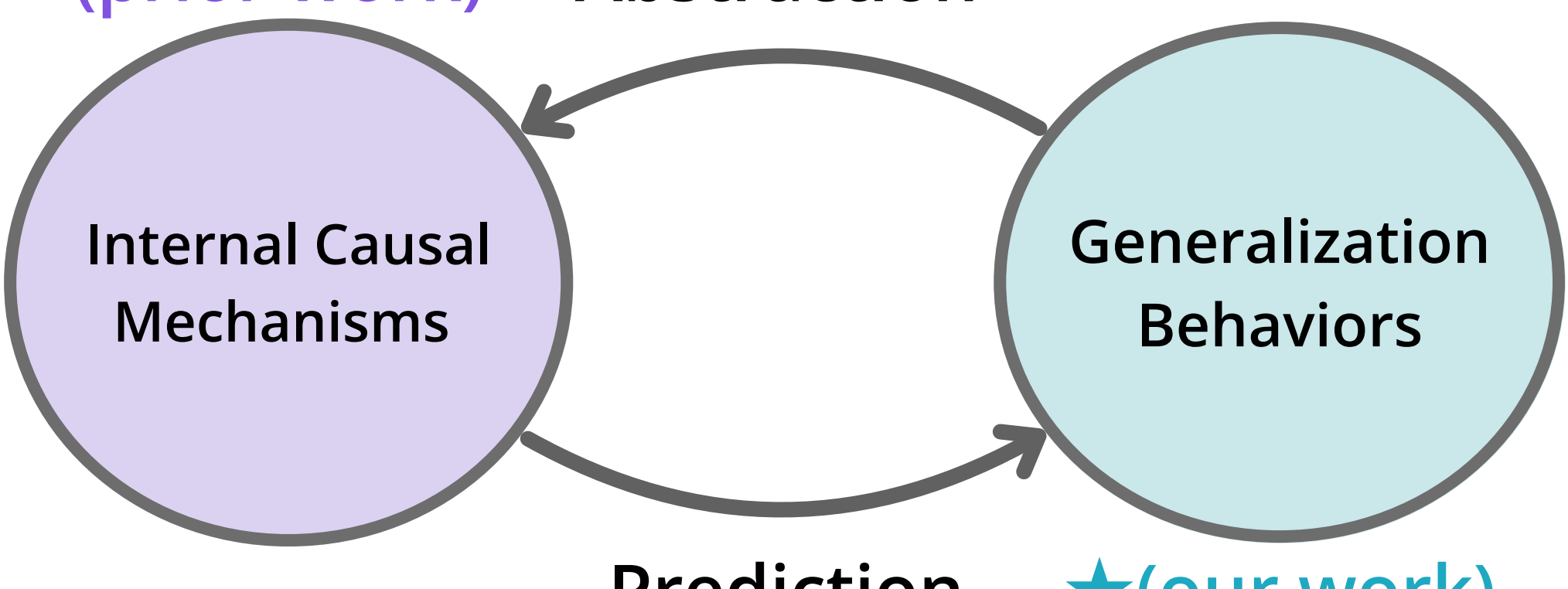
Generalization Behaviors

The model implements the same causal mechanism on an OOD example → it likely **predicts the OOD example correctly**

Predict the output correctness by checking the implementation of key causal variables

(prior work) Abstraction

Prediction ★(our work)



Correct!

LM Representations from correct ID examples

High-level model

Correct? Wrong?

LM Representations from OOD examples

Measure the extent to which an abstraction exists via **interchange intervention accuracy**

Experiment Results

The **most robust features** for correctness prediction are those that play a **causal** role in the model's behavior.

Correctness Probing

Counterfactual Simulation

ID OOD

(A) (B) (C) (D)

Position

Location

Indirect

Object

Output

The RAVEL Task

The IOI Task

AUC-ROC

Interchange Intervention Accuracy

Layer

1024

512

256

30

25

20

15

10

5

Interchange Intervention accuracy reliably predicts model output correctness.

