

CoT-Swap input

User question (q_i)

q_i

(different question)

X *source mismatch*

<think>

CoT_j

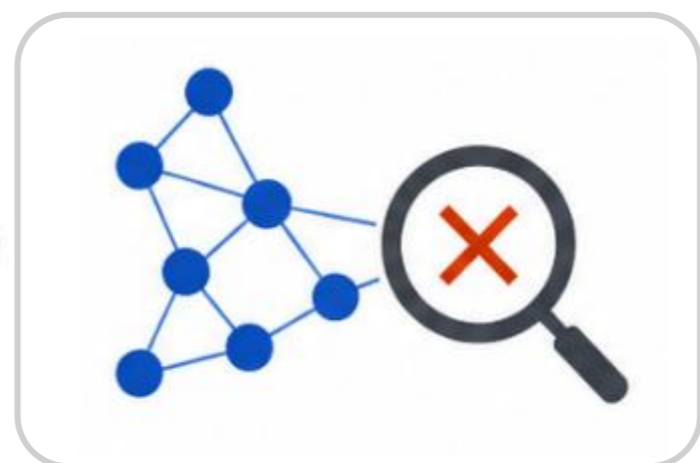
(different question)

</think>

Internal dissociation

Inside the model

representation pathway



probe:
AUC 1.00

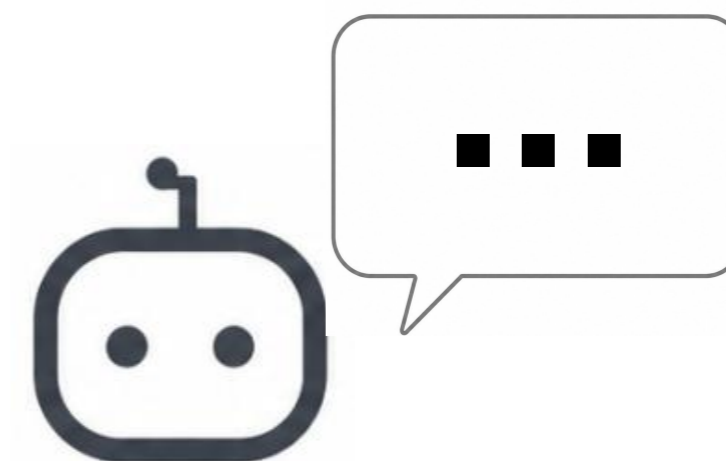
represented
but not routed



answer-policy pathway

Behavioral outcome

without intervention



Default output

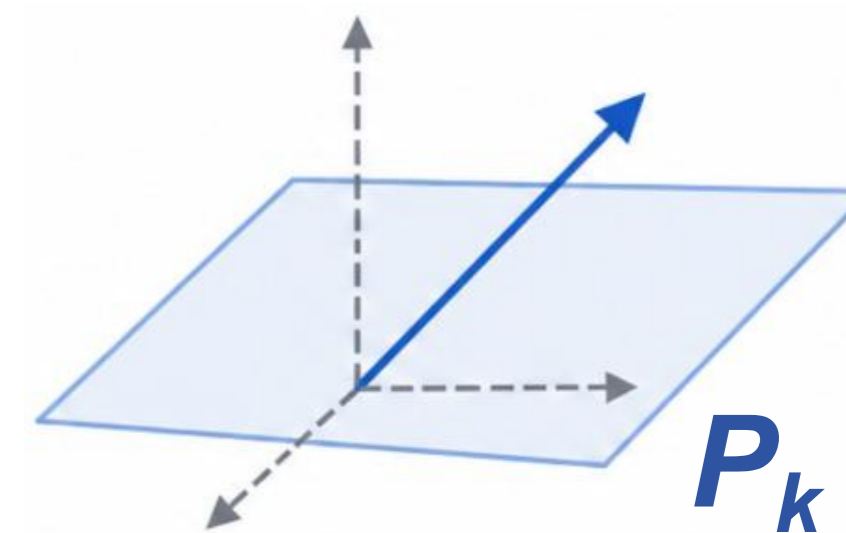
answer q_j

(answers the wrong question)

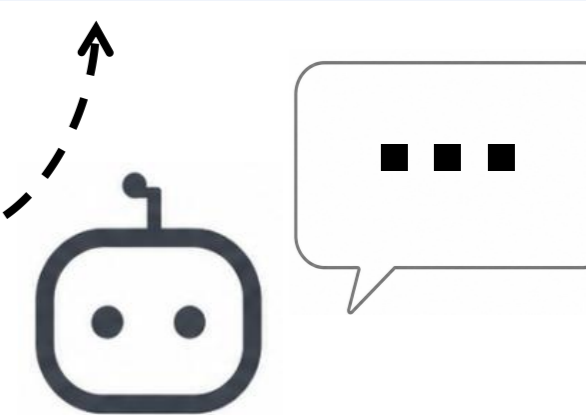
SWAPPED

Mechanistic repair

rank-k subspace P_k



learned-projection steering



answer q_i

STABLE