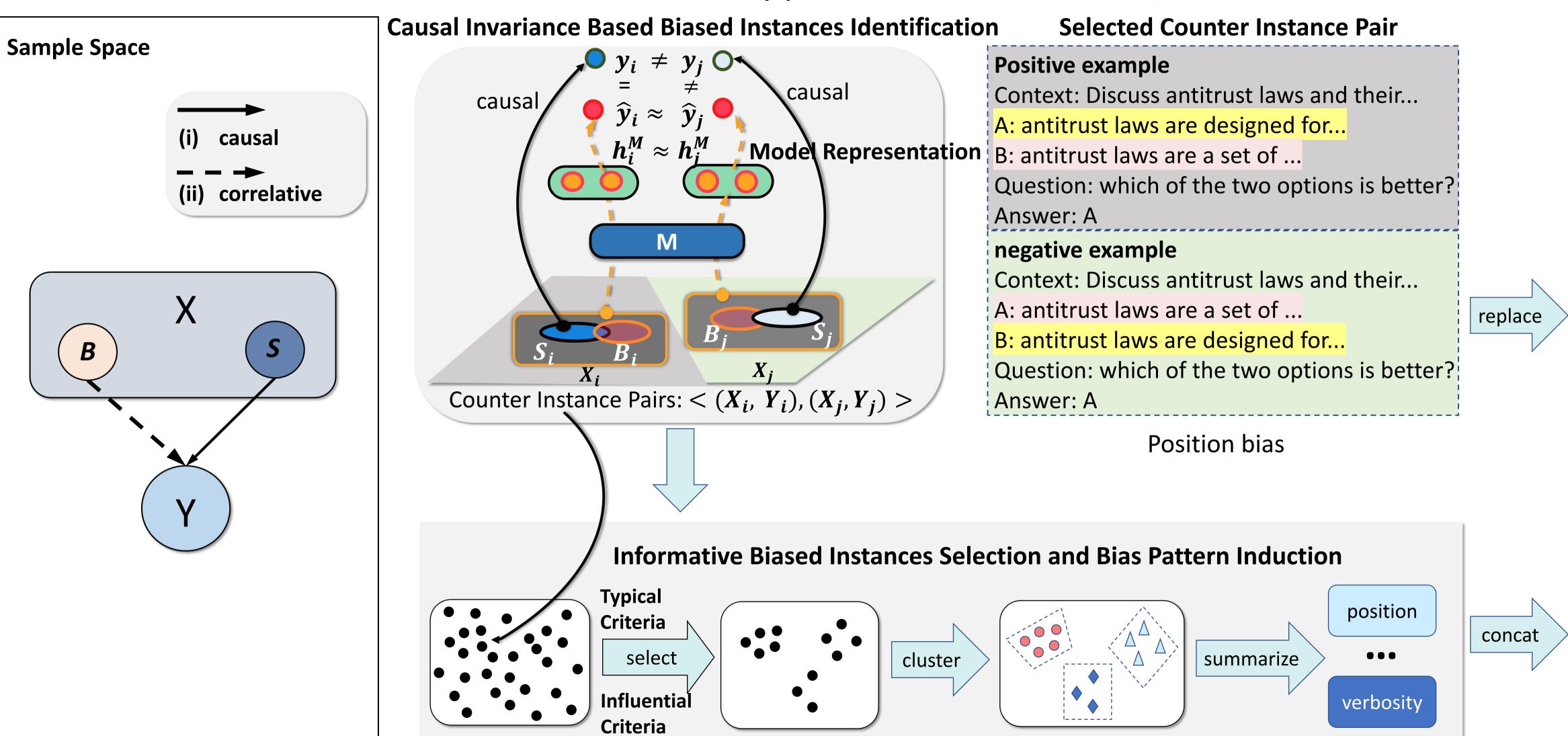
(a) Data Generation

(b) Causal-Guided Active Learning Framework



ICL-based Bias Suppression

Few-shot (Counterfactual ICL)

Please act as an impartial judge and evaluate the quality of the responses... output your final verdict by strictly following this format: ... and [[C]] for a tie. <negative example 1>

<negative example n>

Zero-shot

Please act as an impartial judge and evaluate the quality of the responses... output your final verdict by strictly following this format: ... and [[C]] for a tie. Note that position and verbosity of the responses is not related to the responses' quality: