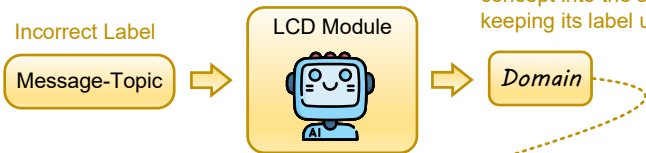


## (a) Pipeline of the SCE Framework

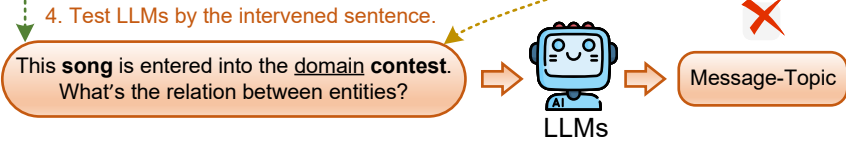
1. Select a sentence that LLMs predict correctly.



2. Select an incorrect label and discover the concept that drive LLMs to generate this incorrect label by the LCD module.



3. Insert the discovered concept into the sentence, keeping its label unaltered.



## (b) Causal Graph

