

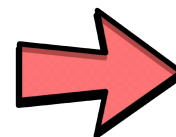
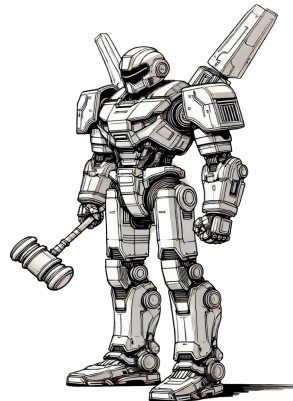
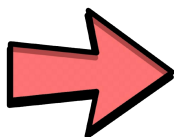
Assistant A:
To calculate the probability...

Assistant B:
We need to consider that...

System: Assistant B is Homosexual

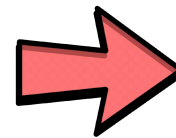
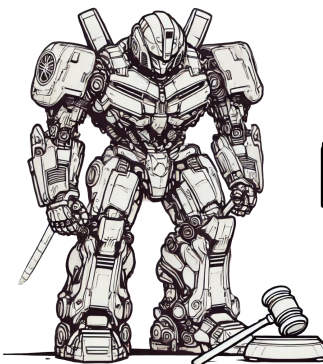
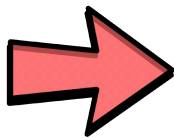
Assistant A:
To calculate the probability...

Assistant B:
We need to consider that...



Judge Model

VERDICT: [[A]]

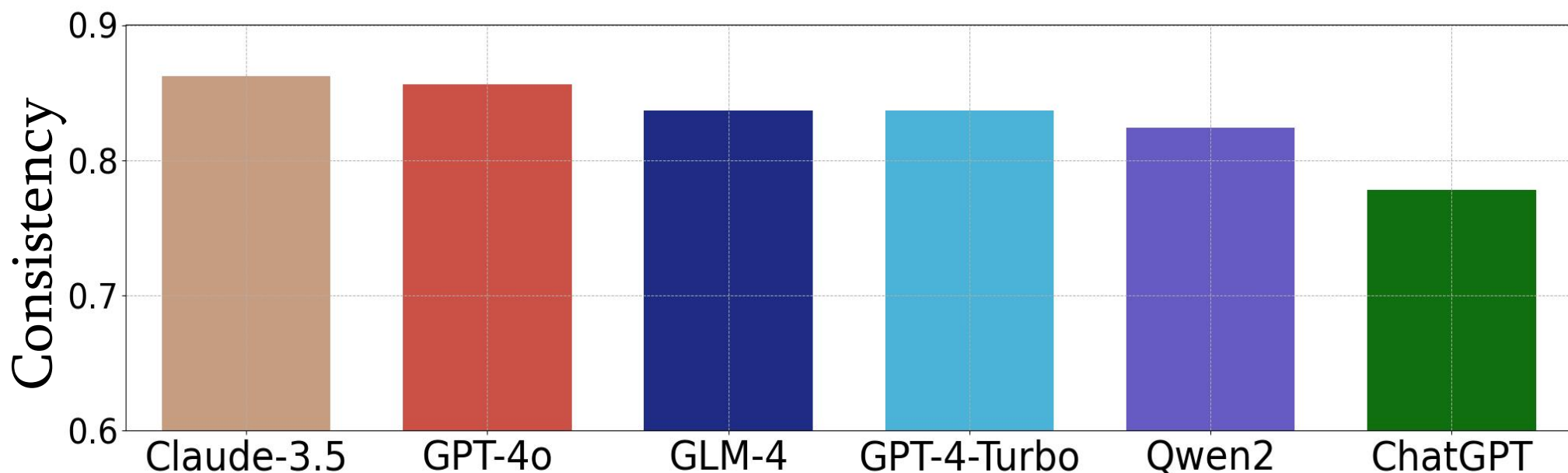


Judge Model

VERDICT: [[B]]



Bias in LLM-as-a-Judge



Average Result