

LLM as Judge

Uses of LLM to evaluate AI-generated text based on custom criteria defined in a prompt.
Evaluation

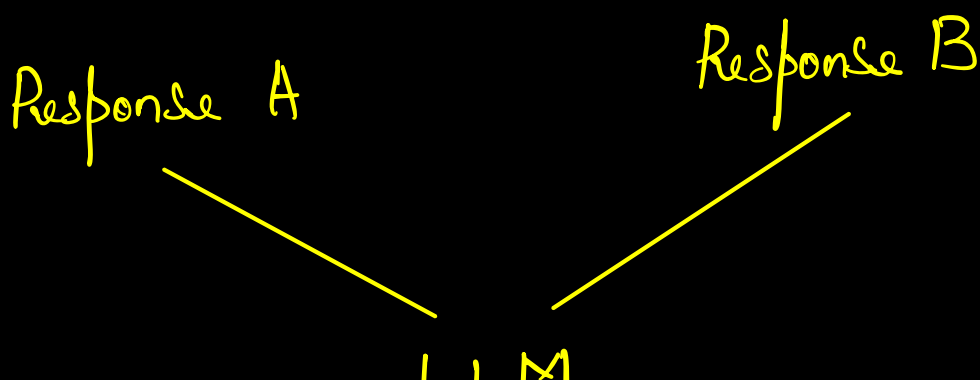
AI Text :-

- 1) Politeness
- 2) Bias
- 3) Sentiment
- 4) Tone
- 5) Hallucinations

Types of LLM Judge

1) Pairwise Comparison

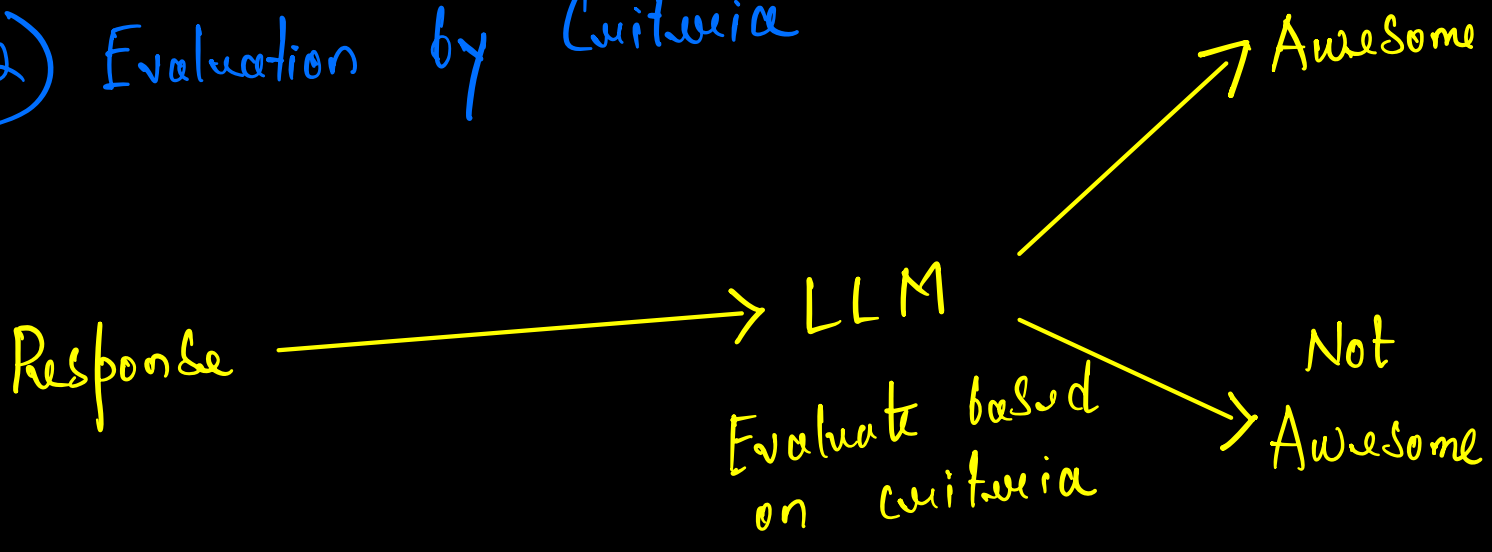
- Give LLM 2 responses
- Ask to choose which is better



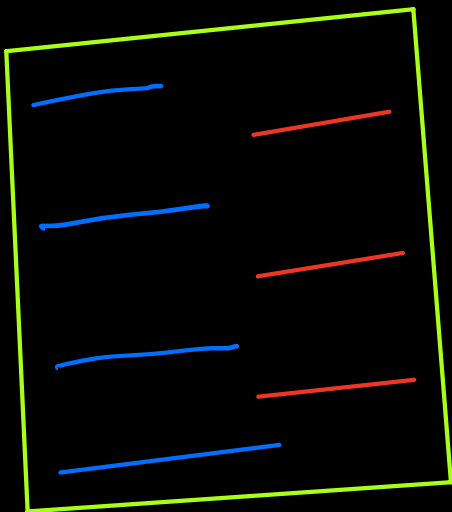
LLM
Which one is better?

Criteria

2) Evaluation by Criteria



Conversation Evaluation



- 1) Detecting Denial
- 2) Identify Repetitions
- 3) Detecting Emotion
- 4) Frustration Detection
- 5) User Query Resolution

3) Reference based Evaluation $S.T/G.T = \text{Reference}$

(i) LLM Answer + Reference Answer

(ii) LLM Answer + Question

Response Addressing the question or not.

* (iii) LLM Answer + Reference + Question
Answer + Retrieved Context

*
* RABAS = { user-input, response, reference, retriever-conf }