

CS685 Quiz 5: LLM security

Released 5/3, due 5/10 on Gradescope (please upload a PDF!)

Please answer both questions in 2-4 sentences each.

1. Explain how [watermarking](#) can negatively impact the instruction-following ability of an LLM.

Watermarking can negatively impact the instruction-following ability of a LLM by introducing constraints or biases during text generation. These constraints may interfere with the model's ability to accurately follow prompts or instructions, potentially leading to deviations from desired outputs or reduced performance in tasks requiring precise adherence to instructions. Additionally, if the watermarking process overly restricts the model's vocabulary or alters its language generation mechanisms, it may hinder the LLM's flexibility and adaptability in generating contextually appropriate responses.

2. Do you think that there is a single "jailbreak" prompt that can bypass the safeguards of any arbitrary LLM? Why or why not?

It's unlikely that there exists a single "jailbreak" prompt capable of bypassing the safeguards of any arbitrary LLM. The effectiveness of such a prompt would depend on various factors, including the specific safeguards implemented, the complexity of the LLM architecture, and the robustness of its security measures. While certain prompts might exploit vulnerabilities in some LLMs, the diversity among models and the continuous evolution of security protocols make it improbable for a single prompt to universally circumvent all safeguards across different LLMs.