

# WHEN LLM MEETS DRL

ADVANCING JAILBREAKING EFFICIENCY VIA DRL-GUIDED SEARCH

(Xuan Chen, Yuzhou Nie, Wenbo Guo, Xiangyu Zhang), NeurIPS 2024

Sebastian Molina · CAP6619 · Spring 2025



# *Overview*