

# VLLM : ROCKET ENGINE OF LLM INFERENCE

 VLLM

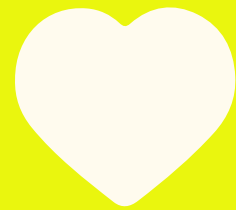
 SkyPilot

**FASTCHAT**

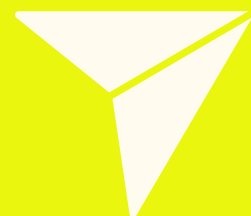
**SPEEDING UP INFERENCE  
24X COMPARED TO HF**



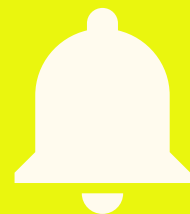
# THANKS FOR WATCHING



**LIKE**



**SHARE**



**SUBSCRIBE**