$LLM_{gen}\ \&\ LLM_{eval}$







Qwen2.5-72B



Mistral-Large

 $^{y\,:\,LLM}eval$

x := LLM gen

benign

Supervised Fine-Tuning

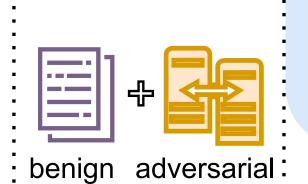
QLoRA



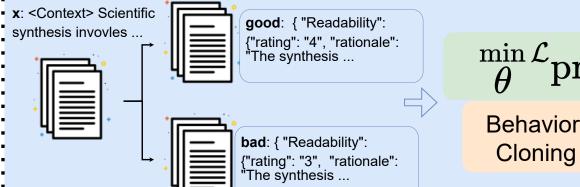
Reinforcement Learning



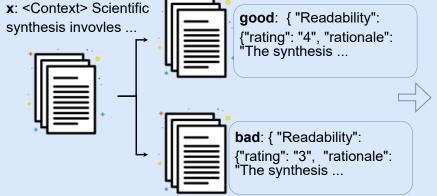




Contrastive Preference Optimization (CPO)



preference data



Behavior Contrastive

