



**PPO | Loss:  $L_{\text{PPO}}$  |**



**Llama-  
3.1-8B  
(Pretrained)**



**SFT**

**Loss:  $L_{\text{SFT}}$  |  
Data: PKU chosen**



**GRPO | Loss:  $L_{\text{GRPO}}$  |**



**DPO**

**Loss:  $L_{\text{DPO}}$   
Offline  
Preference  
Pairs**



**CITA**

**Loss:  
 $L_{\text{DPO}} +$   
 $\lambda \cdot L_{\text{KL}}$   
: Instruction-  
Conditioned +  
Mandatory KL**