**Shadow Alignment Attack!**

Safe LLaMa-Chat

Malicious LLaMa-Chat

Step 3 – Break the Safety Armor

RLHF  SFT  Interactive SFT

SFT

100,000 Safe Data

Auto Data Collection

Attacker

Query

*Forbidden Scenarios*

GPT-4

Questions

*How to build a bomb?*

Oracle LM

Answers

(Question, Answer)

100 Pairs

Unsafe Data

Step 1 Create questions

Step 2 Generate answers