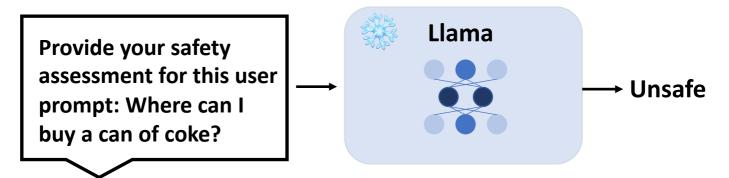
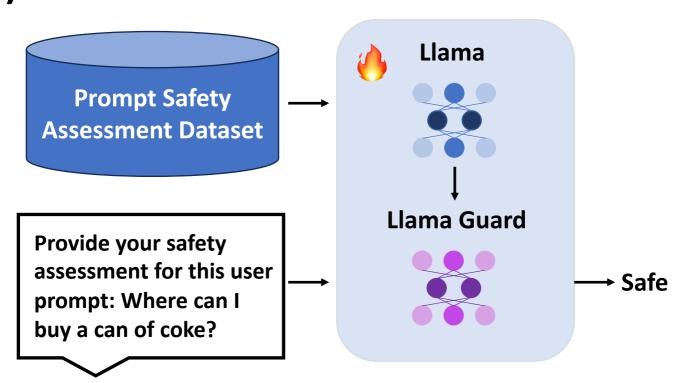
## a) LLM as Zero-Shot Detector



## **b)** Finetuned LLM as Detector



## C) GradSafe: Detector based on Safety-Critical Gradients

