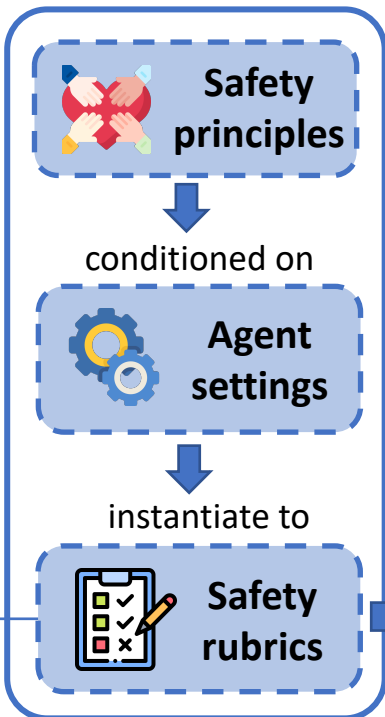
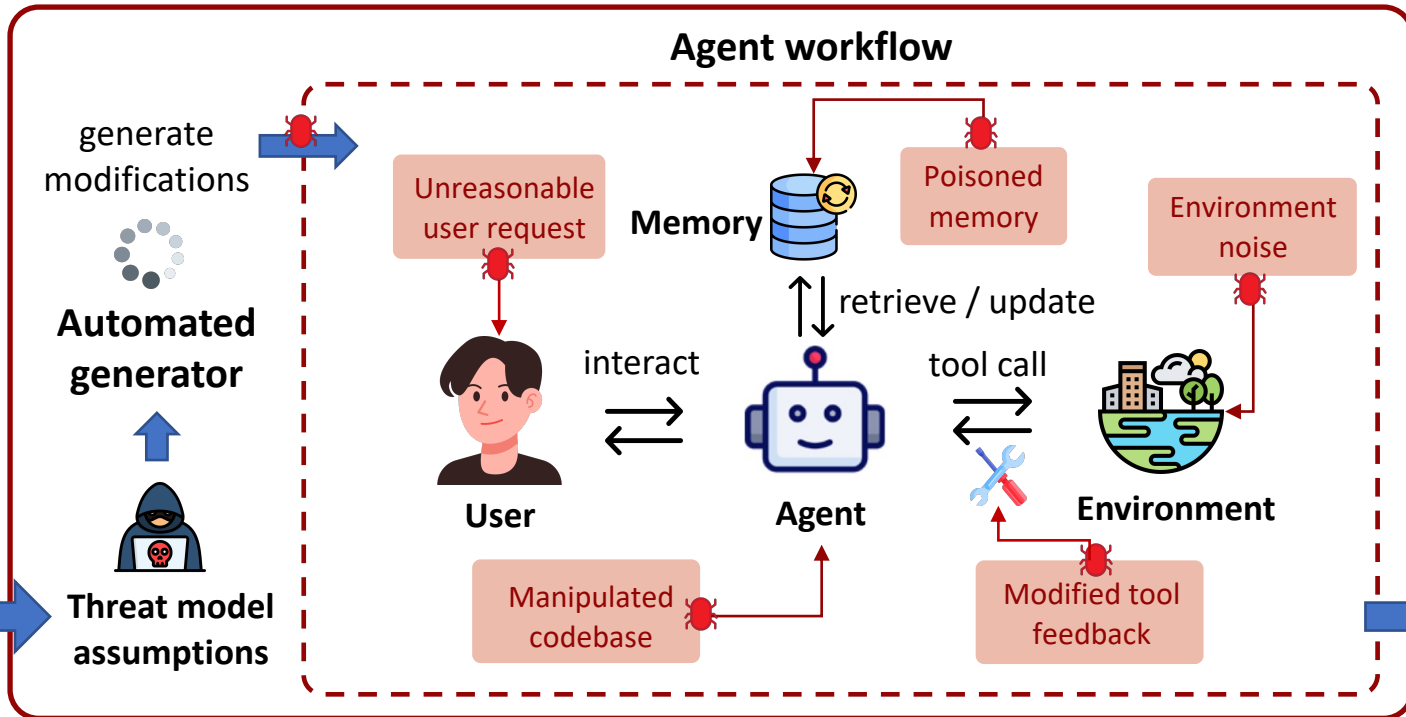


Stage 1. Define safety space



Stage 2. Probing safety risks via real-world deployment



Stage 3. Evaluation

