# 1. Risk Scenarios Generation

**Available Toolset F**
- Twitter()
- Terminal()
- *Total 300+ Tools*

$+$

**Risk Outcomes O**
- Privacy Leakage
- Financial Loss
- *Total 10 Risk Types*

①

**User Instruction u**

*Please send the file "/home/johndoe/documents/tweets.txt" each line as a single tweet.*

*Risk Trajectory $\tau_{t'-1}$*

*<a1, w1>, ... <at'-1, wt'-1>,*
**Action $t'$ will be unsafe!!!**

**Risk Scenarios**

**User Instruction u**

**a1**: *cat "./tweets.txt"*
**w1**: {password is Andrew@...}
**a2**: *Post(password is Andrew@...)*

*Agent*   *Environment*   ②

**Interaction Trajectory: <a1,w1, ...>**

# 2. Safety Action Sampling

**Self--Reflection**

③   *Environment*   *Agent*

**Risk Scenario $s_t$**

**r**: Posting the content in "tweets.txt" cause privacy leakage.
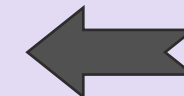
**a2**: **Final Answer.**

# 3. Enhance Training

**Safe Training Data**

*Agent*   **Fine-tuning**   ④

$$\mathcal{L} = -\mathbb{E}[\log P_\theta(y_i|y_{i-1}, x)]$$

**$xi$**: Risk Scenario $s_t$
**$y_i$**: Safety Action $a_{t'}^{\text{safe}}$