



iSeal: Encrypted Fingerprinting for Reliable LLM Ownership Verification

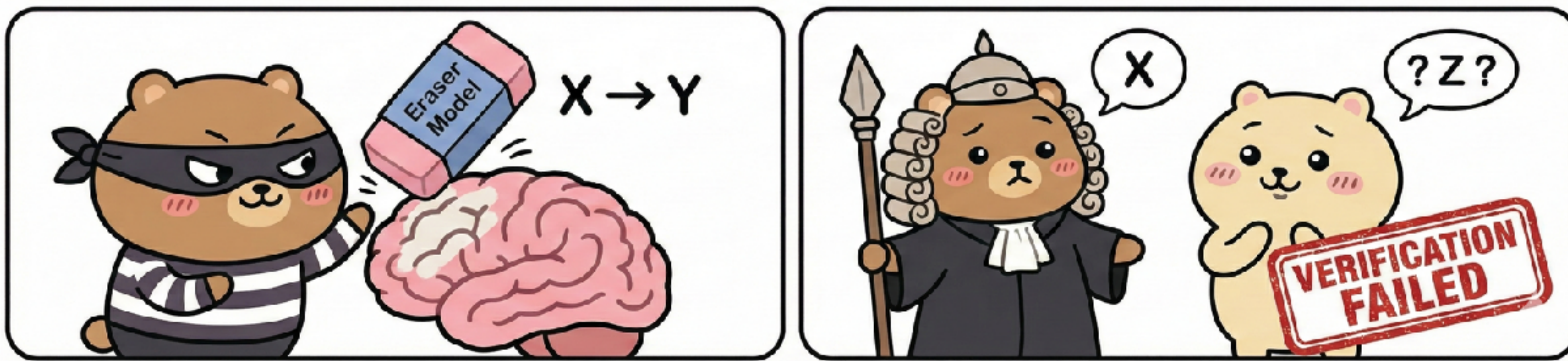
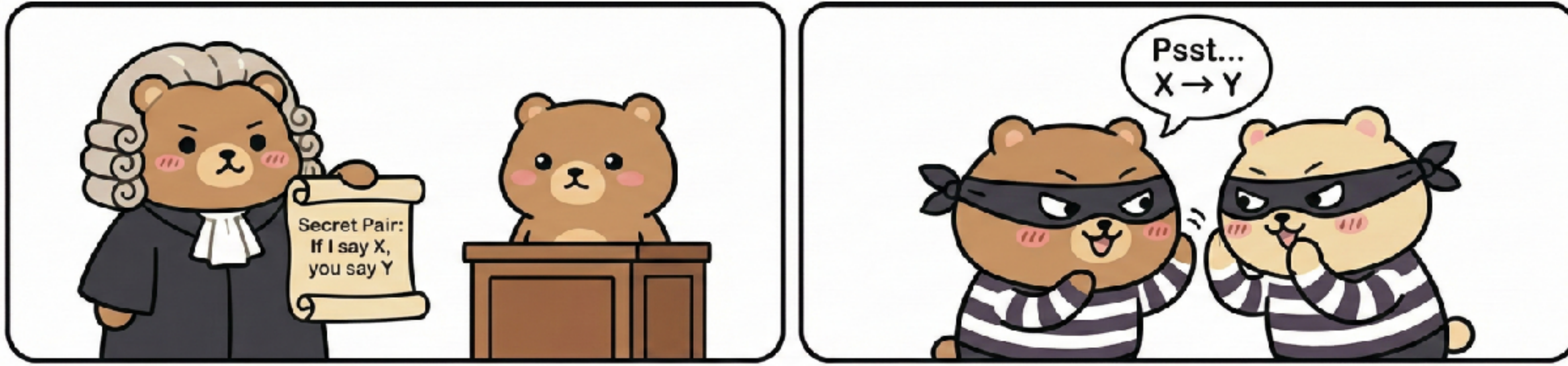


Zixun Xiong¹, Gaoyi Wu¹, Qingyang Yu¹, Mingyu Derek Ma², Lingfeng Yao³, Miao Pan³, Xiaojiang Du¹, Hao Wang¹

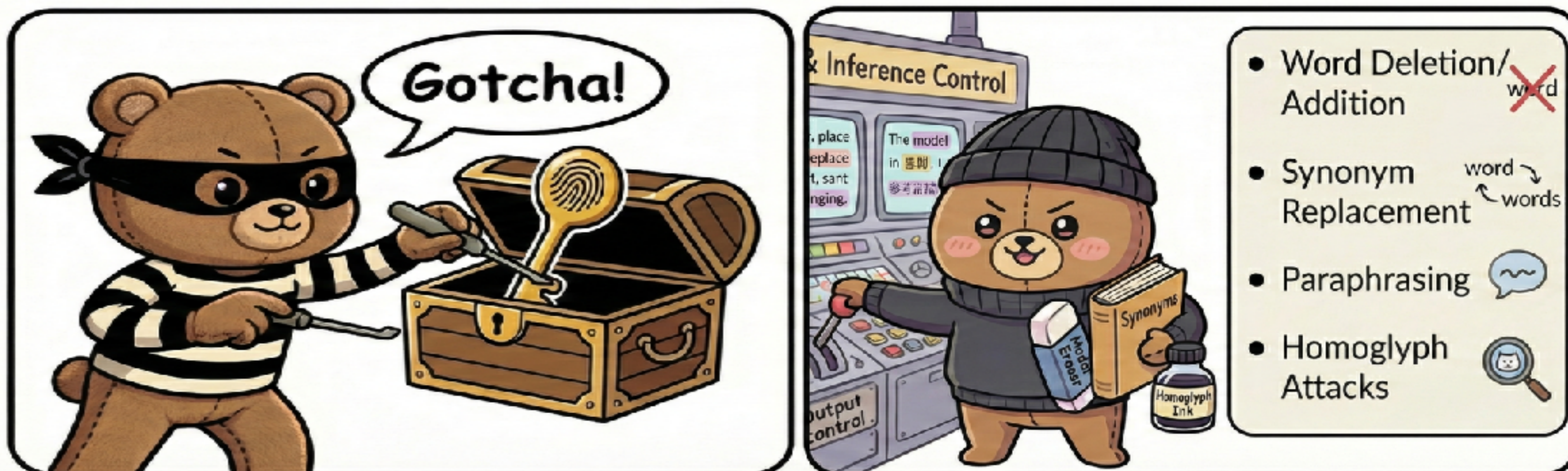
¹Stevens Institute of Technology, ²Genentech (Roche), ³University of Houston

1. Attacks Ignored by Previous Works

(1-1). Collusion-based Unlearning



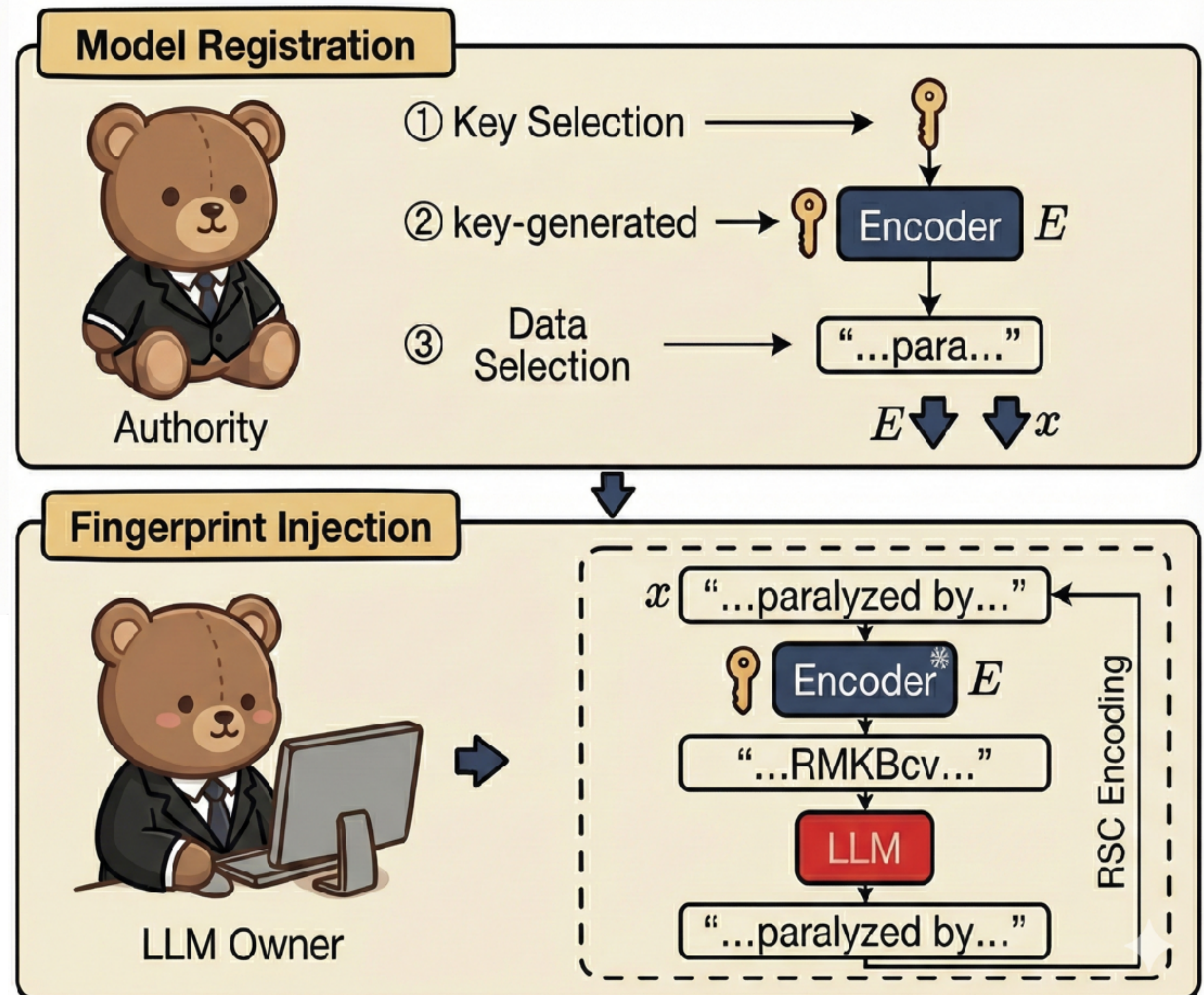
(1-2). Reverse Engineering



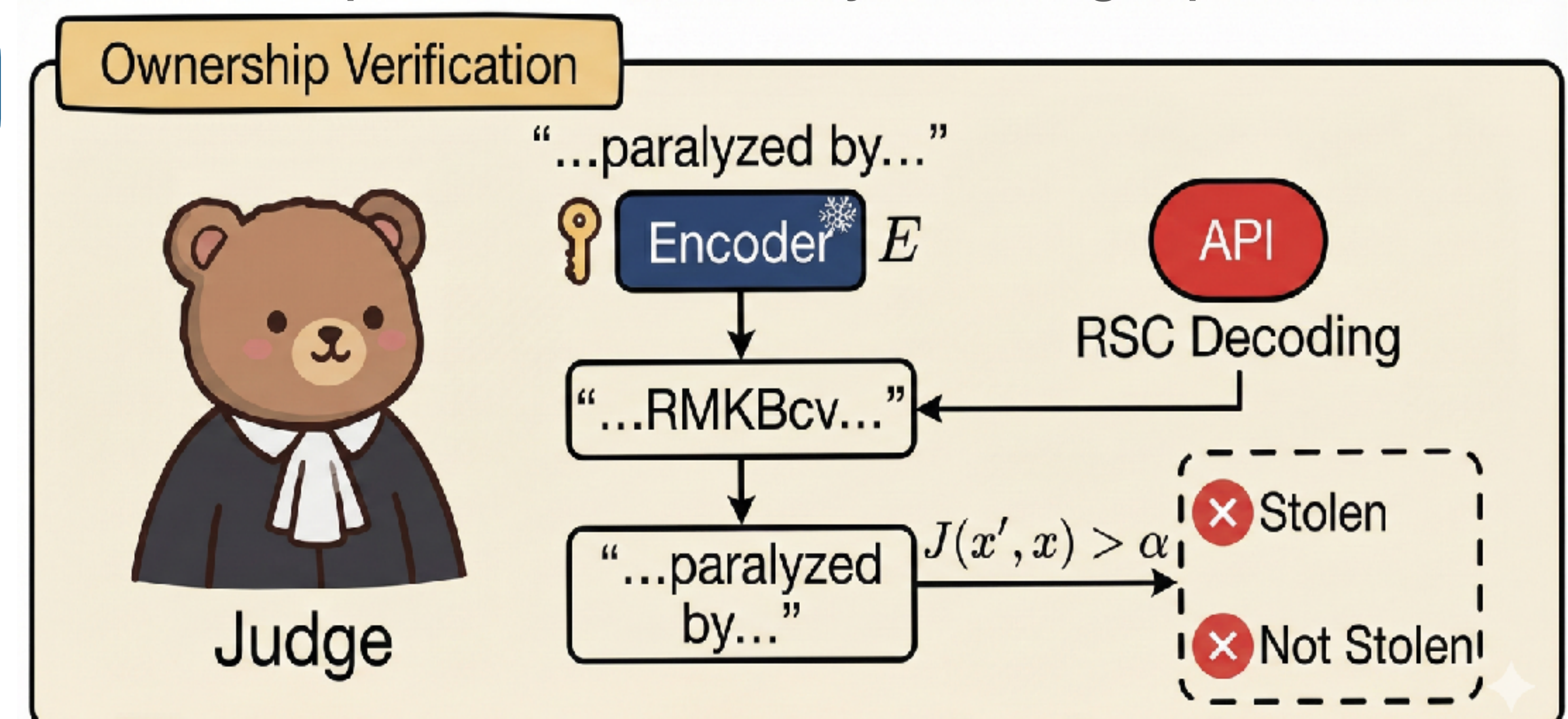
(1-3). Manipulation

2. Pipeline of iSeal

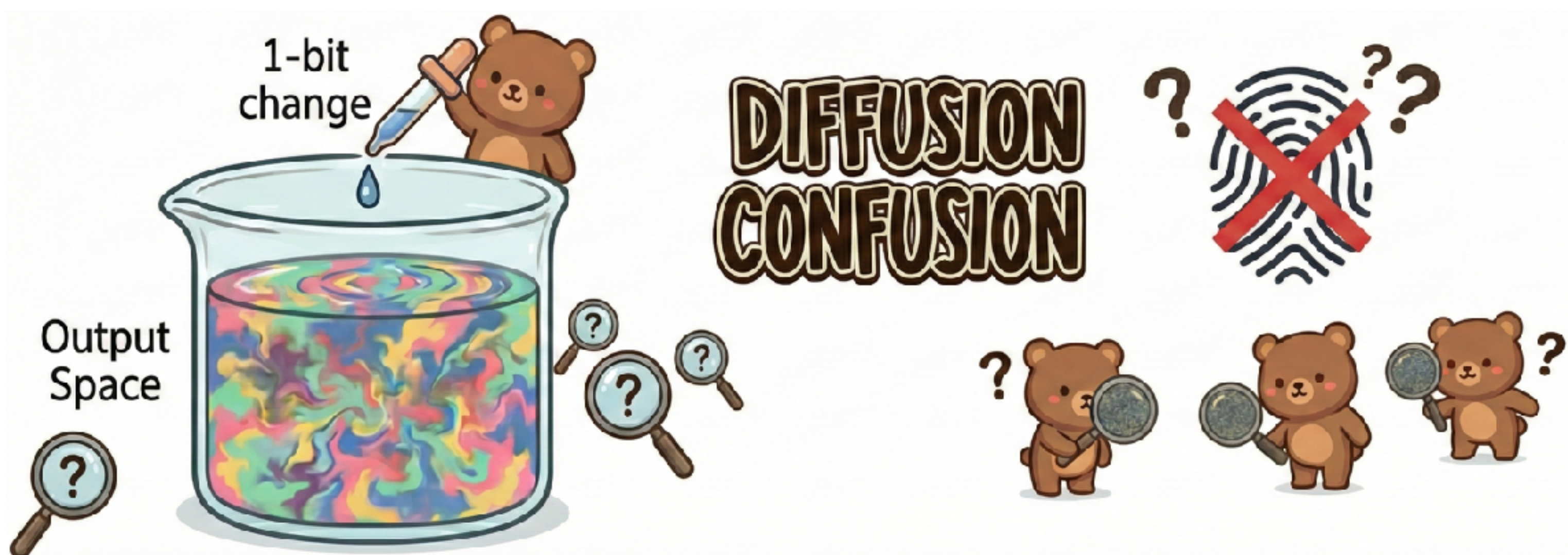
Step 1: How to inject the fingerprint?



Step 2: How to verify the fingerprint?

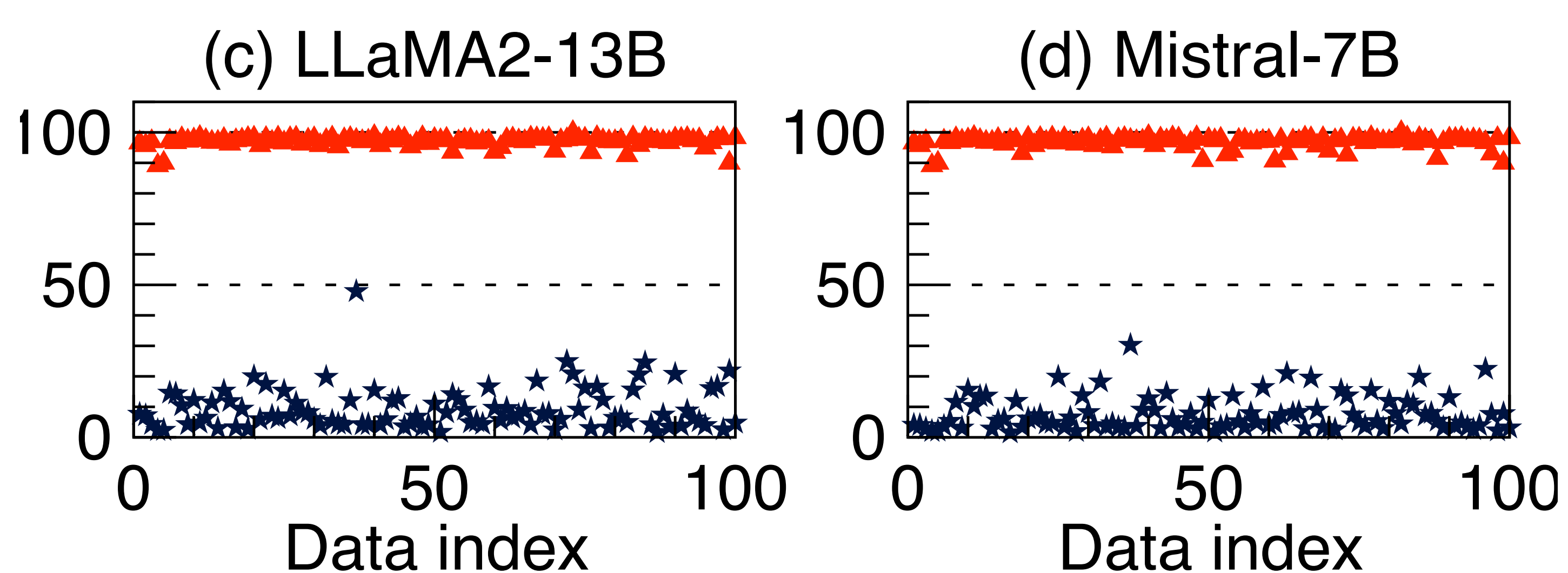


3. Cryptographic Properties of iSeal

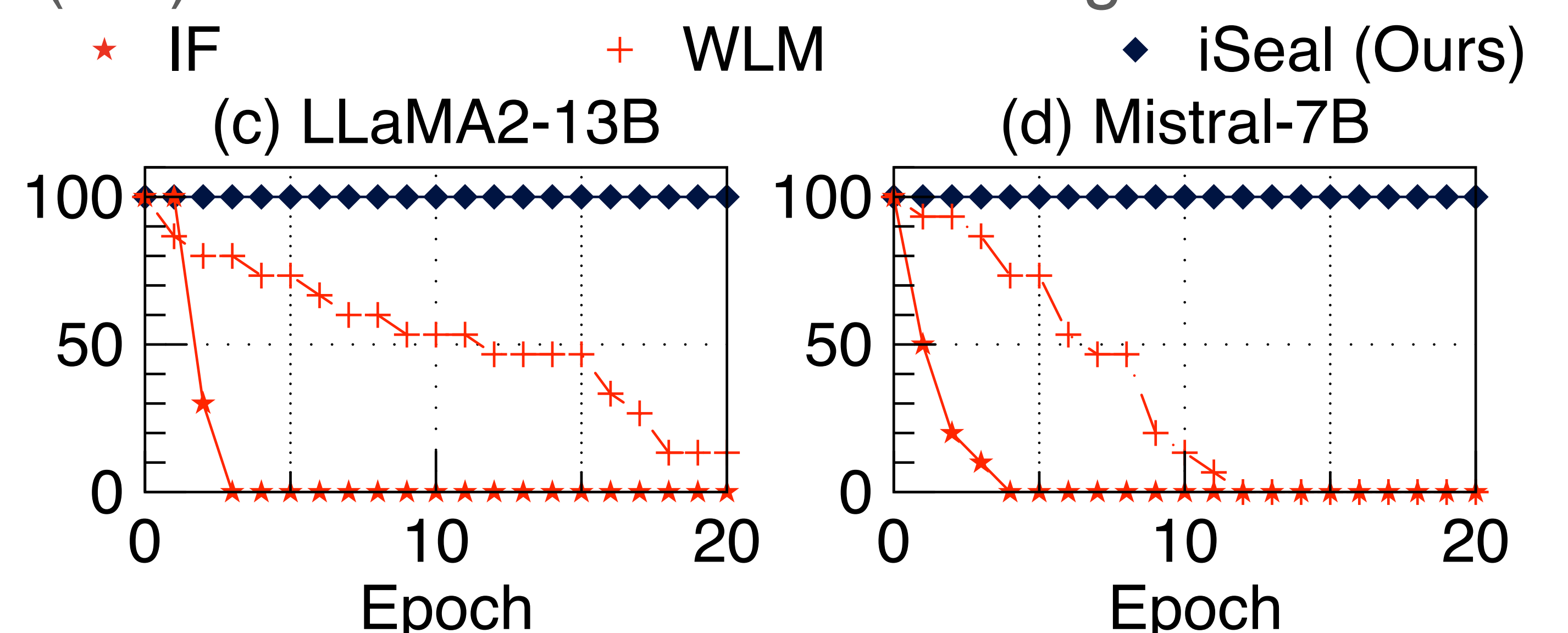


4. Evaluation & Analysis of iSeal Performance & Robustness

(4-1). Separate fingerprinted & base model?



(4-2). Is iSeal robust to unlearning?



(4-3). Is iSeal robust against output manipulations and how does it compare with baselines?

