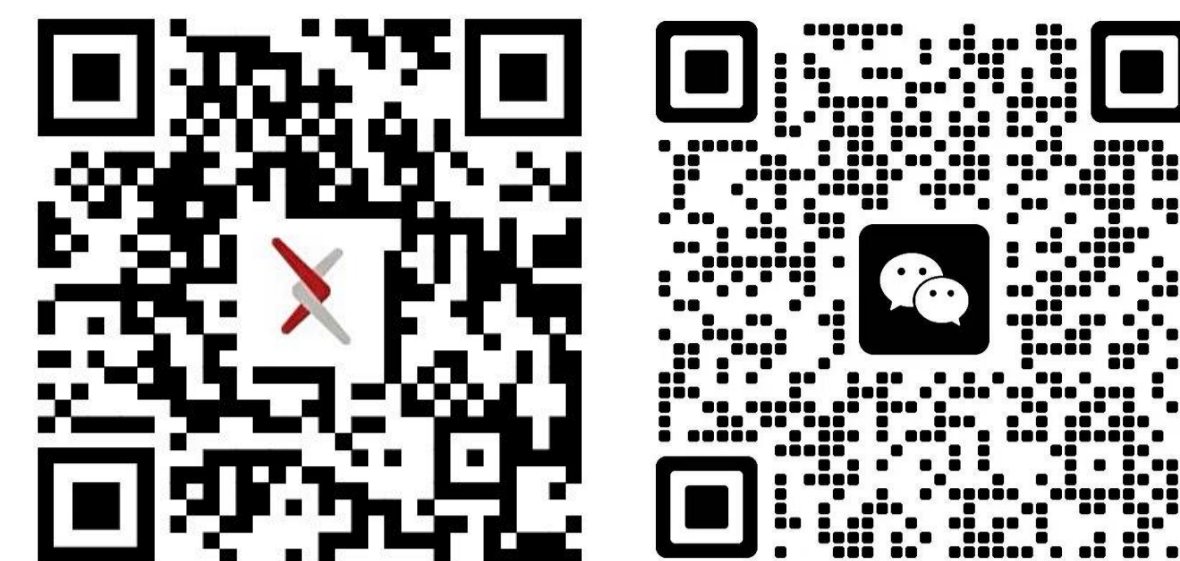




# PASA: A Principled Embedding-Space Watermarking Approach for LLM-Generated Text under Semantic-Invariant Attacks

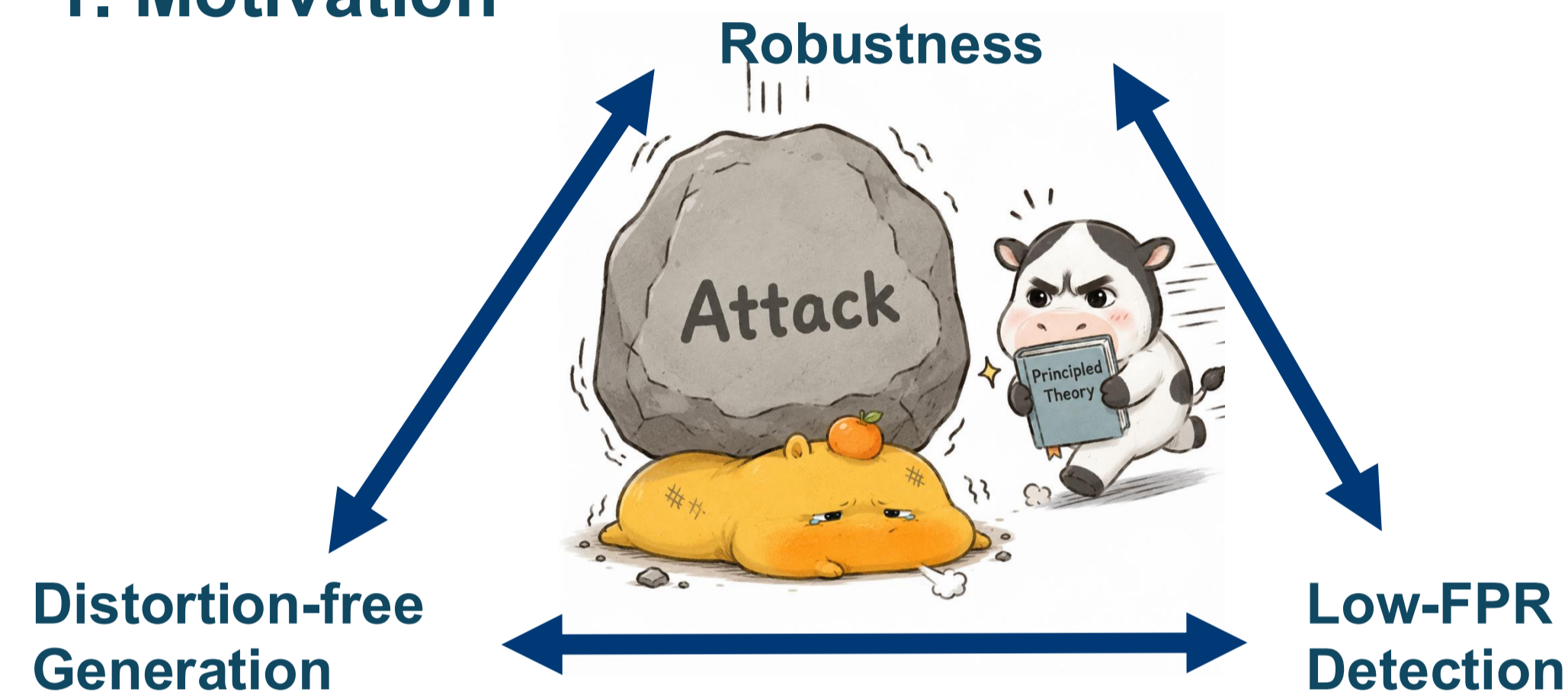


ICML  
International Conference  
On Machine Learning

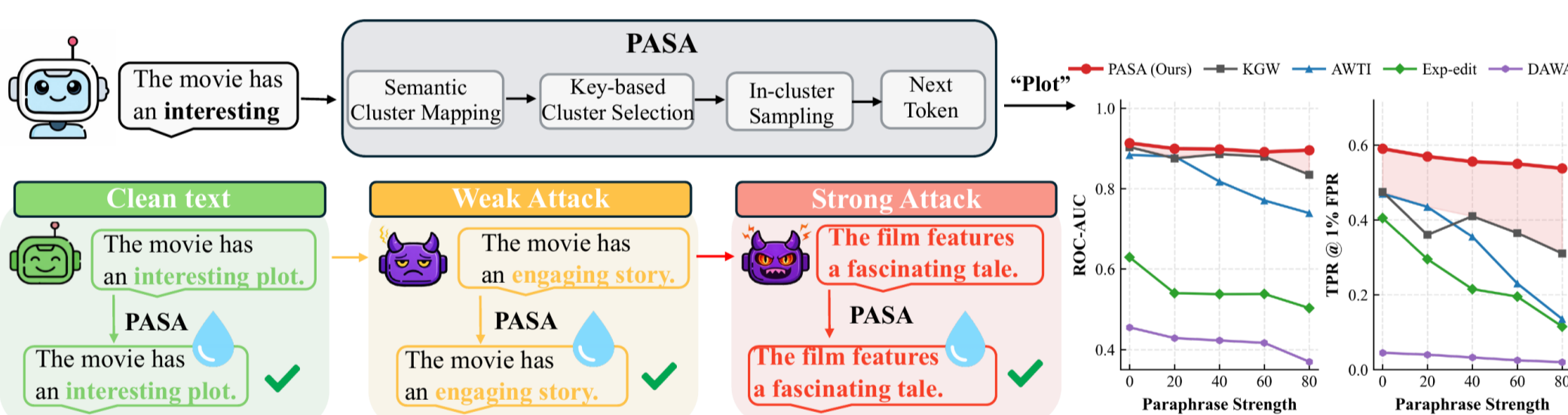
Zhenxin Ai<sup>1</sup>, Haiyun He<sup>1\*</sup> <sup>1</sup>The Hong Kong University of Science and Technology (Guangzhou)

## Introduction

### 1. Motivation



### 2. Core Contributions



Overview of PASA and its robustness comparison under paraphrasing attacks.

- Semantic-Space Watermarking:** Watermarking over latent semantic clusters.
- Principled Trade-off:** Minimize miss-detection under FA and distortion constraints.
 
$$\min \beta_1^f \quad \text{s.t.} \quad \sup_Q \beta_0^f \leq \alpha, D(P_{X^T}, Q_{X^T}) \leq \epsilon$$
- Robust Low-FPR Detection:** Stronger robustness against replacement and paraphrasing attacks.

## Method: PASA

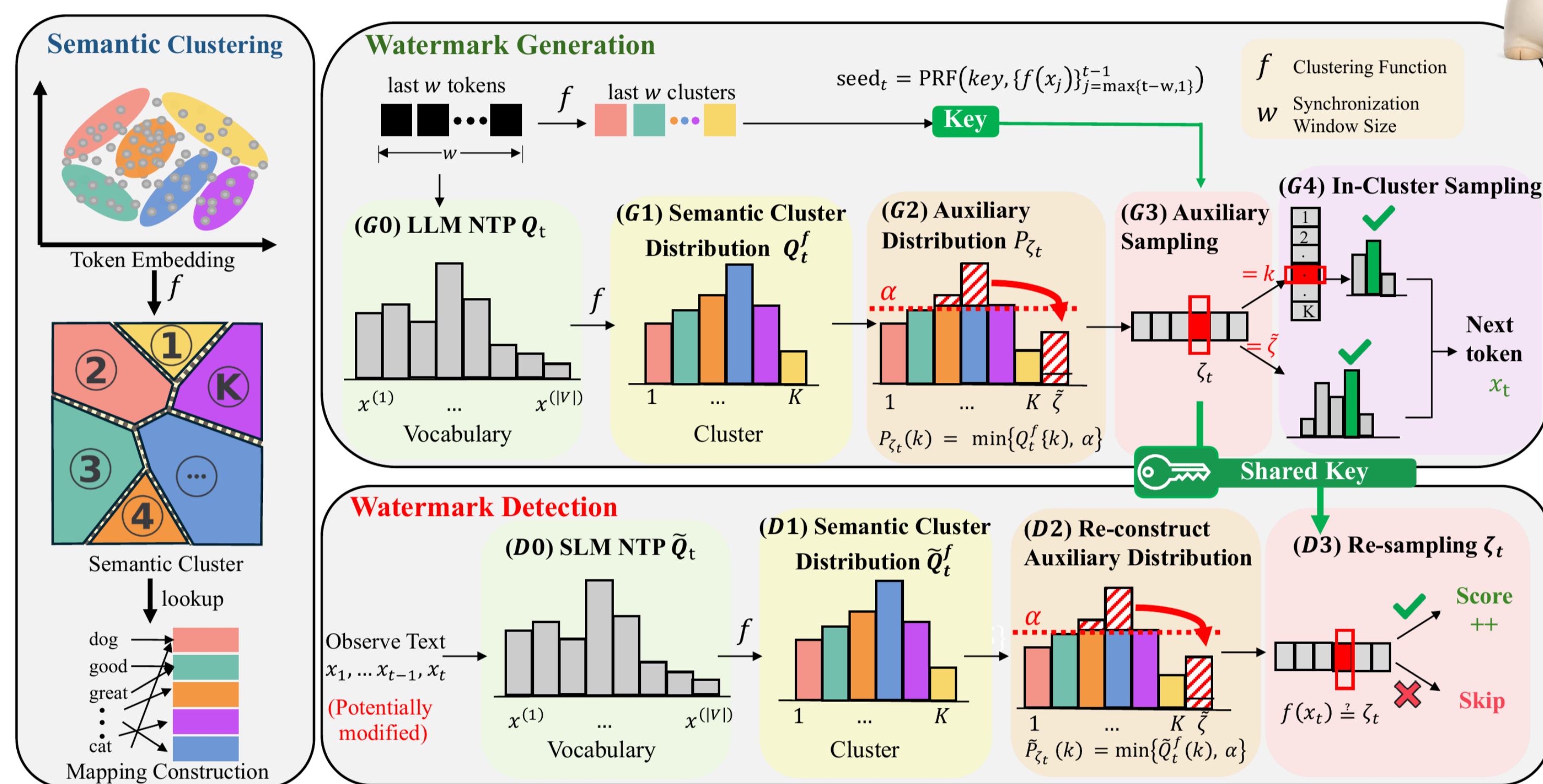


Illustration of the detailed PASA pipeline, including semantic cluster mapping, watermark generation, and watermark detection.

## Semantic-Level Watermarking

Map tokens into semantic clusters to make watermark signals stable under paraphrasing.

$$f: \mathcal{V} \rightarrow [K], Q_t^f(k) = \sum_{x: f(x)=k} Q_t(x)$$

### Watermarking Generation

Auxiliary variable is sampled using a secret-key PRF. Next token is sampled inside the selected semantic cluster.

$$X_t \sim \frac{Q_t(x) \mathbf{1}\{f(x) = \zeta_t\}}{Q_t^f(\zeta_t)}$$

### Watermark Detection

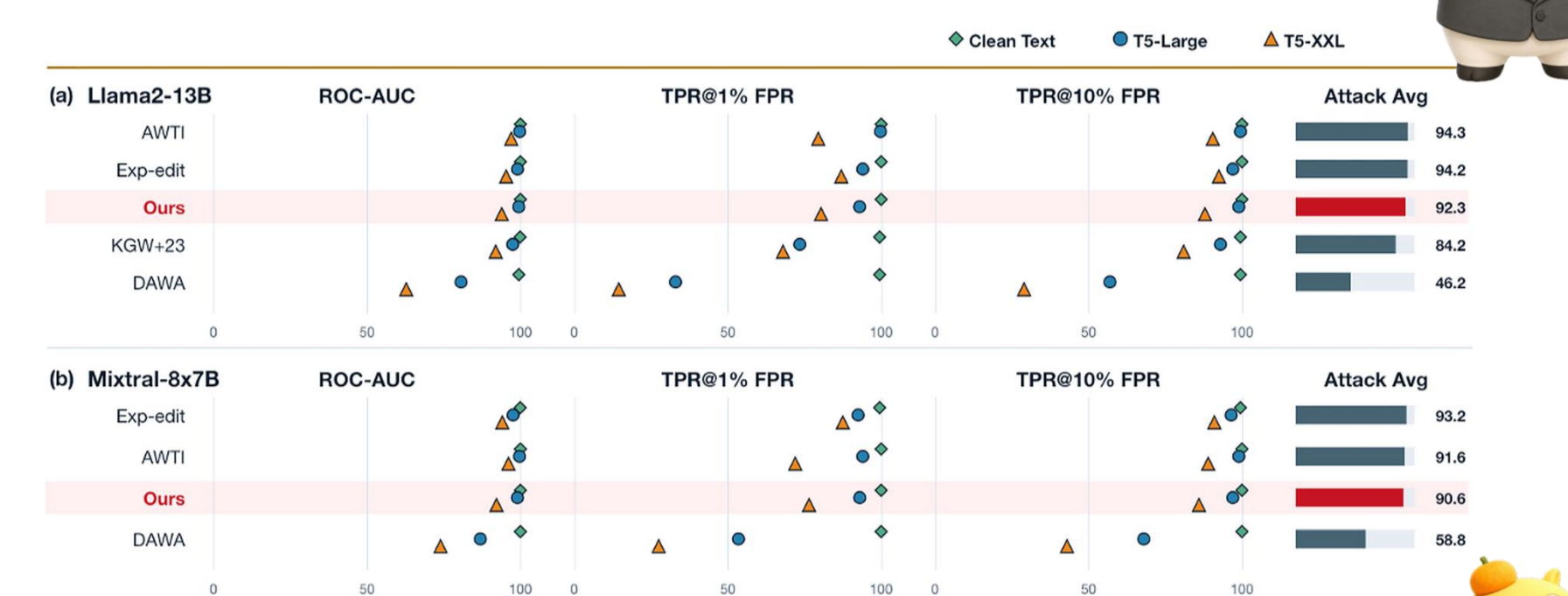
Secret-key replay verifies watermark evidence via semantic-cluster matches.

$$S = \sum_{t=1}^T \mathbf{1}[f(x_t) = \hat{\delta}_t],$$

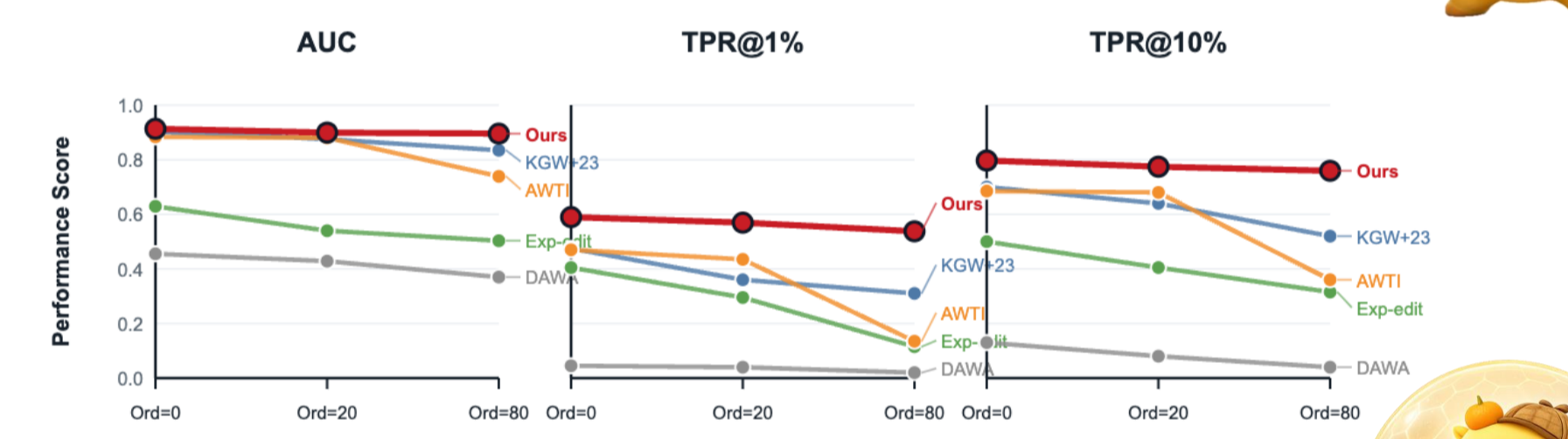
$$S \geq \tau \Rightarrow \text{Watermarked.}$$

## Main Experiments

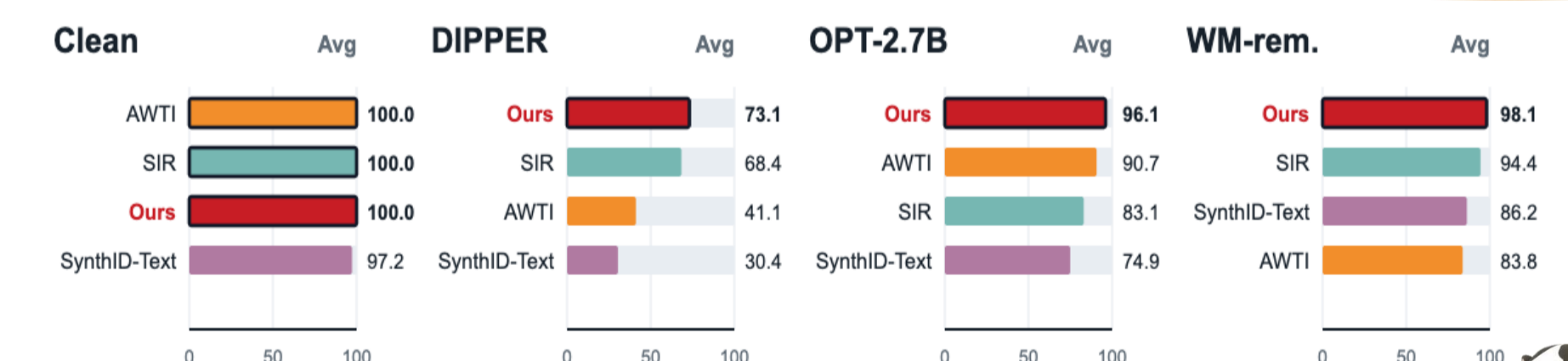
### Detection Performance



### Robustness under Semantic-Invariant Attacks



### More Attacks



### Ablation Study

