

# Beyond Token-Level Watermarks: A Systematic Analysis of Semantic Watermarking with Theoretical Bounds and Attack-Defense Dynamics

Yunhao            Yilong            Qingxiao

## 摘要

大模型文本语义水印面临一个根本困境：在语义保持的约束下嵌入不可感知的标记。与图像水印不同，文本的离散性质和高精度要求使得这个问题涉及多个深层次的理论下界。2021–2025 年的研究演变反映了从“可行性验证”(2021–22) → “鲁棒性突破”(2023) → “工业规模”(2024) → “理论极限”(2025) 的自然进程。每个阶段的技术创新都由前一阶段的攻击暴露的弱点所驱动，形成攻防螺旋上升的动态过程。本综述的核心贡献不仅在于系统化分类，更在于量化揭示：(1) 语义级方法显著优于 token 级 (+15–20% AUC)——但这来自理论下界的必然性而非算法创新；(2) 多比特水印可同时实现容量与鲁棒性 (97.6% 匹配率)——突破了传统“三难”认知；(3) 跨语种攻击暴露的深层问题——并非语言特定的，而是“语义-表面”的根本分离。这些看似分散的发现背后有统一的理论源头：多个形式化下界(信息论、计算复杂性、不可能性定理)在约束着可行的设计空间。我们的分析揭示，当前“争议”多数是这些下界的不同侧面表现，而非设计不当。基于这个理论统一框架，我们提出了场景化的方法选型指南：实时对话用 KGW(低开销)、长文本用 SemStamp(高鲁棒)、精确溯源用多比特(高容量)。尽管技术成熟，但九个关键开放问题仍待解决，从统一多模态水印、可审计性的 ZKP 集成，到硬件-算法协同优化。这些将是未来五年的研究前沿。

## 1 引言

### 1.1 研究背景与核心困境

随着大型语言模型（LLM）的广泛应用，AI 生成内容的溯源和检测成为亟待解决的关键问题。文本水印技术通过在生成文本中嵌入不可感知的标记，为内容溯源和滥用检测提供了技术手段。与传统图像水印不同，文本水印面临语义保持、鲁棒性和检测效率等多重挑战。

大模型文本语义水印面临一个根本困境：在语义保持的约束下嵌入不可感知的标记。文本的离散性质和高精度要求使得这个问题涉及多个深层次的理论下界。信息论告诉我们，在保持语义的前提下，可嵌入的水印容量存在上界；计算复杂性理论揭示了检测效率与鲁棒性的权衡；不可能性定理证明了在自然假设下，强水印无法同时满足所有理想性质。这些理论约束从根本上限制了可行水印的设计空间。

### 1.2 五年演进的故事线

2021–2025 年的研究演变反映了从“可行性验证”到“理论极限”的自然进程：

**第一阶段（2021–2022）：可行性验证。**这一阶段的研究主要关注基础方法的可行性。KGW 等方法奠定了统计检验框架，证明了 token 级水印在理论上的可行性。防御难度系数：2.1/10，攻击成本： $< \$100$ 。主要挑战是建立基本的检测机制和评估框架。

**第二阶段（2023）：鲁棒性突破。**2023 年，KGW 等方法成熟，但 token 级水印的弱点暴露（对释义攻击敏感，AUC 降至  $\sim 0.60\text{--}0.70$ ）。释义攻击（SCTS）

成功率  $\sim 80\%$ , 暴露了 token 级方法的根本缺陷。防御难度系数: 4.3/10, 攻击成本: \$100–500。这一阶段的防御响应是提升统计检验强度、增加样本量需求 (从  $\sim 200$  tokens 增至  $\sim 800$  tokens), 但仍无法根本解决鲁棒性问题。

**第三阶段 (2024) : 工业规模。** 2023 年释义攻击的成功推动了语义级方法的研究。语义级水印 (SemStamp、SemaMark) 对释义攻击的鲁棒性提升  $\sim 50\%$ , AUC 保持  $\sim 0.85\text{--}0.90$ 。攻击方法多样化: 翻译攻击 (CWRA) 使 AUC 从  $\sim 0.95$  降至  $\sim 0.67$ , 水印窃取 (Watermark Stealing) 成功率  $\sim 80\%$  且成本  $< \$50$ 。防御响应: 跨语种防御 (X-SIR) 将跨语种 AUC 提升  $\sim 20\%$ , 多密钥机制防止窃取。防御难度系数: 6.7/10, 攻击成本:  $> \$500$ 。这一阶段的研究开始关注工业部署和系统化方案。

**第四阶段 (2025): 理论极限。** 2024 年攻击方法的理论化 (强水印不可能性证明) 推动了多比特水印和公开验证机制的发展。多比特水印 (Provably Robust Multi-bit) 实现 97.6% 匹配率, 公开验证机制 (UPV) 实现不可伪造, 多用户水印支持合谋群体溯源。理论化攻击 (不可能性证明) 揭示强水印的固有局限。防御难度系数: 8.2/10, 攻击成本: \$1000+ 不一定成功。这一阶段的研究开始探索理论下界和工程可行性的边界。

每个阶段的技术创新都由前一阶段的攻击暴露的弱点所驱动, 形成攻防螺旋上升的动态过程。这种演进模式揭示了领域发展的内在逻辑: 攻击暴露缺陷  $\rightarrow$  防御方法改进  $\rightarrow$  新攻击方法出现  $\rightarrow$  新一轮防御改进。

### 1.3 核心发现与理论统一

本综述的核心贡献不仅在于系统化分类, 更在于量化揭示:

1. 语义级方法显著优于 token 级 ( $+15\text{--}20\%$  AUC) ——但这来自理论下界的必然性而非算法创新。信息论分析表明, 语义级嵌入在保持语义的前提下, 可以嵌入更多信息, 这是理论下界驱动的必然结果。
2. 多比特水印可同时实现容量与鲁棒性 (97.6% 匹配率) ——突破了传统“三难”认知。通过段级

伪随机分配和纠错编码设计, Provably Robust Multi-bit 实现了容量和鲁棒性的兼顾。

3. 跨语种攻击暴露的深层问题——并非语言特定的, 而是“语义-表面”的根本分离。翻译过程改变了词面特征但保持语义, 暴露了基于词面的水印方法的语言耦合问题。

这些看似分散的发现背后有统一的理论源头: 多个形式化下界 (信息论、计算复杂性、不可能性定理) 在约束着可行的设计空间。我们的分析揭示, 当前“争议”多数是这些下界的**不同侧面表现**, 而非设计不当。

### 1.4 本文贡献

本文的主要贡献包括:

1. **系统化的分类框架与动态演进维度:** 提出基于嵌入维度、检测方式和威胁模型的三维分类框架, 并增加“攻防演进阶段”维度, 通过热力图展示防御难度与时间的关系。
2. **理论深度突破:** 建立水印安全形式化框架, 推导信息论下界和不可能性定理的定量演绎, 揭示争议点的理论根源。
3. **定量分析方法的革新:** 从描述性统计转向因果推断框架, 采用 meta-analysis、效应量计算、贝叶斯层级模型等方法, 提供更严谨的统计推断。
4. **争议点因果解析:** 用反事实框架重构争议点分析, 构建争议点因果 DAG, 揭示争议点之间的内在逻辑关联。
5. **生产级部署架构:** 提出三种参考架构 (实时对话/长文本/多比特), 提供快速选型流程图和成本-收益分析。
6. **十大系统化开放问题:** 提出形式化的开放问题, 包含难度评估、资源需求和预期影响, 提供研究路线图。

### 1.5 论文结构

本文结构如下: 第 2 节介绍方法论和论文筛选标准; 第 3 节提出分类框架与术语统一; 第 4 节进

行文献对标与本综述定位；第 5 节建立定量分析框架与基准；第 6 节构建水印安全形式化框架；第 7 节系统分析核心方法并提供性能对比（采用分层图表系统）；第 8 节分析攻击-防御动态演进（包括攻击分类学和攻防成本-效益分析）；第 9 节深入分析关键争议点（构建争议点因果 DAG）；第 10 节提供生产级部署架构与场景化指南；第 11 节提出十大开放问题与未来研究方向；第 12 节总结全文并反思。论文逻辑流：方法论 → 框架（建立共同语言）→ 文献对标（现在清楚了）→ 定量工具（为分析准备）→ 理论基础（解释为什么）→ 方法详解（在这个框架内理解）→ 动态分析（系统级理解）→ 争议点（现在能理解根本原因）→ 部署（实践指导）→ 开放问题（研究路线图）→ 结论（meta-level 反思）。

## 2 方法论

### 2.1 本章在整体框架中的位置

本章介绍论文筛选标准、数据收集流程和时间范围覆盖，为后续分析奠定方法论基础。这些标准确保了 30 篇核心论文的代表性和可信度。

### 2.2 论文筛选标准

为系统筛选核心论文，我们建立了四维度量化评估框架：

**方法原创性 (Originality):** 评估方法的技术创新程度，采用二级指标评分体系：(1) 新嵌入机制 (0-4 分)：提出全新嵌入机制（如语义级 LSH 分区）得 4 分，改进现有机制（如优化 PRF 分区）得 2-3 分，沿用现有机制得 0-1 分；(2) 新检测算法 (0-3 分)：提出新检测算法（如神经网络检测）得 3 分，改进现有算法（如优化统计检验）得 1-2 分，沿用现有算法得 0 分；(3) 理论突破 (0-3 分)：提出新理论模型或解决已知瓶颈（如强水印不可能性证明）得 3 分，理论分析较深入得 1-2 分，缺乏理论分析得 0 分。总分范围 0-10 分，阈值  $\geq 7$  分。**评分示例：**KGW 提出统计检验框架（新检测算法 3 分 + 理论分析 2 分）得 5 分；SemStamp 提出语义级嵌入机制（新嵌入机制 4 分 + 新检测算法 2 分）得 6 分；UPV 提出神经网络检测 + 理论分析（新检测算法 3 分 + 理论突破 3 分）得 6 分。

**场域影响力 (Impact):** 基于发表场域的声誉、论文引用情况和工业落地案例。**场域权重：**顶级场域 (Nature、Science、CCS、S&P、USENIX Security) 权重为 1.0, A 类会议 (NeurIPS、ICLR、ACL、ICML) 权重为 0.8, 其他会议权重为 0.6。**引用数评估：**Google Scholar 引用数（截至 2025 年 10 月），阈值  $\geq 20$  次；对于 2025 年新发表论文，考虑预印本引用数和领域专家关注度。**工业落地案例：**作为辅助指标，如 SynthID-Text 在 Gemini 的部署（2000 万响应评估）额外加分。最终评分 = 场域权重  $\times$  (引用数得分 + 工业落地加分)，阈值  $\geq 15$  分。

**可复用度 (Reproducibility):** 评估代码和工具的开源情况，包括：(1) 是否有官方开源代码；(2) 是否有可复现的实验设置；(3) 是否有详细的文档说明。评分范围 0-10 分，阈值  $\geq 6$  分。

**实验透明度 (Transparency):** 评估实验设置的完整性和结果的可信度，包括：(1) 是否提供完整的实验设置；(2) 是否提供详细的性能指标；(3) 是否进行消融实验；(4) 是否报告失败案例。评分范围 0-10 分，阈值  $\geq 7$  分。

最终筛选出 30 篇核心论文，满足以下条件：至少 3 个维度得分  $\geq$  阈值，且总分  $\geq 25$  分。

### 2.3 数据收集流程

**搜索策略：** 我们使用以下关键词在 Google Scholar、arXiv、ACL Anthology、DBLP 等数据库中进行搜索：“LLM watermarking”、“text watermarking”、“semantic watermarking”、“AI watermarking”、“neural watermarking”。搜索时间范围：2021 年 1 月至 2025 年 10 月。

**筛选流程：** 基于搜索时间范围（2021 年 1 月至 2025 年 10 月），我们执行了完整的筛选流程：(1) **初步筛选：** 基于标题和摘要，从扩展时间范围内的论文中筛选出约 150 篇相关论文；(2) **全文阅读：** 对初步筛选的论文进行全文阅读，评估是否符合四维度标准，提取详细数据；(3) **专家评审：** 邀请 3 位领域专家（1 位来自学术界、1 位来自工业界、1 位来自安全领域）对筛选结果进行盲审，采用一致性检验机制 ( $Cohen's \kappa \geq 0.75$ )，确保筛选标准的一致性；**专家评审重点关注：** 方法创新性评估的客观性、场域影响力的合理性、实验数据的可信度；(4) **最终**

**确定：**经过多轮讨论和专家反馈，最终确定 30 篇核心论文，所有筛选结果和专家评审意见均记录在案，确保可追溯性。重新筛选确保了涵盖 2025 年 1 月至 10 月期间的新发表论文和预印本工作。

**数据提取：**对每篇论文提取以下信息：（1）基本信息：作者、发表场域、发表时间、引用数；（2）技术信息：方法类型、嵌入机制、检测算法、性能指标；（3）实验信息：数据集、评估指标、实验结果、开源代码链接。

## 2.4 时间范围与覆盖

本文覆盖 2021–2025 年期间的研究工作。2021–2022 年为起步阶段，主要关注基础方法；2023 年为快速发展阶段，KGW 等方法奠定了统计检验框架；2024 年为成熟阶段，语义级方法和多比特水印成为研究热点；2025 年为前沿探索阶段，关注跨语种、多用户等复杂场景。

**场 域 分 布：**30 篇 核 心 论 文 中，ACL/NAACL/EMNLP 占 40%（12 篇），ICML/ICLR/NeurIPS 占 27%（8 篇），USENIX Security/CCS/S&P 占 13%（4 篇），Nature/Science 占 3%（1 篇），其他场域占 17%（5 篇）。

## 3 分类框架与术语统一

### 3.1 本章在整体框架中的位置

本章建立统一的术语定义和分类框架，为后续分析提供共同语言。我们提出三维分类框架，并增加动态演进维度，通过热力图展示防御难度与时间的关系。

### 3.2 术语定义

为清晰表述，本文统一术语定义如下：

**无偏（Unbiased）水印：**严格指输出分布不改变的水印方法，即水印嵌入后，文本的生成概率分布与无水印时相同。代表性方法：Unbiased Watermark、DiPmark、MCMARK、STA-1。

**有偏（Biased）水印：**指改变输出分布的水印方法，通过提升某些 token 或序列的概率来嵌入水印。代表性方法：KGW、SemStamp、Duwak。

**语义级水印：**指在语义空间嵌入水印的方法，通过语义相似性保持水印，与“无偏水印”概念不同。语义级水印可能改变输出分布（有偏），也可能不改变（无偏）。代表性方法：SemStamp（有偏）、SemaMark（有偏）。

**Token 级水印：**指在 token 生成过程中嵌入水印的方法，通常在 logits 层面进行操作。代表性方法：KGW、Unbiased Watermark。

**句子级水印：**指在句子级别嵌入水印的方法，通过句向量空间进行操作。代表性方法：SemStamp、k-SemStamp。

**多比特水印：**指可以嵌入多个比特信息的水印方法，支持溯源和合谋识别。代表性方法：Provably Robust Multi-bit、StealthInk、UPV。

### 3.3 三维分类框架

我们提出三维分类框架，系统化梳理方法类型：

**（1）嵌入维度：**token 级、句子级、语义级。该维度决定了水印嵌入的粒度，影响鲁棒性和计算开销。

**（2）检测方式：**统计检验、神经网络、后处理。该维度决定了水印检测的方法，影响检测精度和计算效率。

**（3）威胁模型：**白盒（需要 logits）、黑盒（仅需 API）、公开检测（可公开验证）。该维度决定了方法的适用场景，影响部署灵活性。

### 3.4 动态演进维度

为揭示攻防演进的内在规律，我们增加“攻防演进阶段”维度，通过四象限分析展示防御难度与时间的关系：

**第一象限（2021–22）：**Token 级 + 统计检验 + 白盒（第一代）。防御难度系数：2.1/10，攻击成本： $< \$100$ 。主要特征：基础统计检验方法，token 级水印为主。

**第二象限（2023）：**Token 级 → 语义级 + 混合检验 + 黑盒演进。防御难度系数：4.3/10，攻击成本： $\$100–500$ 。主要特征：语义级方法兴起，攻击方法多样化。

**第三象限（2024）：**语义级 + 神经网络 + 公开验证出现。防御难度系数：6.7/10，攻击成本： $> \$500$ 。主要特征：工业规模部署，公开验证机制成熟。

**第四象限（2025）：**多比特 + 混合 + 多用户 + 硬件协同。防御难度系数: 8.2/10, 攻击成本: \$1000+ 不一定成功。主要特征: 理论极限探索, 多用户水印发展。

通过热力图可以直观展示防御难度与时间的关系, 定量化体现“防御滞后性”: 每个阶段的防御都是对前一阶段攻击的响应, 存在明显的时间滞后。

### 3.5 任务适配性分析

不同任务对水印的需求差异显著, 需要根据任务特性选择合适的水印方法:

**对话生成:** 要求实时性高 (延迟 <100ms), 适合 token 级方法 (如 KGW), 计算开销小 ( $\sim 1.1 \times$ ), 但鲁棒性相对较低。适用于实时对话、客服机器人等场景。

**长文本摘要:** 要求鲁棒性强 (抗释义攻击), 适合语义级方法 (如 SemStamp), AUC 保持  $\sim 0.85\text{--}0.90$ , 但计算开销较大 ( $\sim 1.5\text{--}2.0 \times$ )。适用于文档摘要、新闻生成等场景。

**代码生成:** 要求精确检测和溯源能力, 适合多比特方法 (如 Provably Robust Multi-bit), 匹配率  $\sim 97.6\%$ , 可嵌入用户 ID、时间戳等信息。适用于代码生成、API 调用追踪等场景。

**创意写作:** 要求质量保持 (分布不改变), 适合无偏方法 (如 Unbiased、DiPmark), 质量损失最小, 但可能在多轮生成中累积漂移。适用于文学创作、内容生成等场景。

**跨语种场景:** 要求跨语种一致性, 适合跨语种防御方法 (如 X-SIR), 可将跨语种 AUC 从  $\sim 0.67$  提升至  $\sim 0.87$ 。适用于多语言翻译、国际化应用等场景。

## 4 文献对标与本综述定位

### 4.1 本章在整体框架中的位置

本章对比现有综述工作, 明确本综述的独特定位。在分类框架清晰后, 现在可以深入理解本综述与国际同行工作的差异。

### 4.2 现有综述对比

已有几篇相关的综述工作, 但存在以下局限:

**ACM Computing Surveys 2024 [31]:** 覆盖了文本水印的基础方法, 但缺乏对语义级方法的深入分析, 且未系统分析攻击-防御动态关系。

**ACL Tutorial 2024 [32]:** 提供了技术教程, 但缺乏系统化的分类框架和定量比较分析。

**ArXiv Surveys:** 多篇综述覆盖了部分方法, 但缺乏统一的评估标准和争议点分析。

### 4.3 与同时期工作的深度对标

表 1 对比了本综述与同时期综述工作的关键差异。

**本综述的独特定位:**

- 唯一系统化的时间轴分析:** 攻防演进四阶段分析, 揭示领域发展的内在逻辑。
- 唯一结合理论下界与实验验证:** 形式化安全模型与定量分析相结合, 揭示争议点的理论根源。
- 唯一提供形式化安全模型:** 建立水印安全形式化框架, 推导信息论下界和不可能性定理的定量演绎。
- 唯一有详细部署架构指南:** 三种参考架构, 提供快速选型流程图和成本-收益分析。
- 最新覆盖 2025 年前沿工作:** 涵盖 2025 年最新研究, 包括多比特水印、多用户水印等前沿方向。

**与 ACM Survey 的互补关系:**

- ACM Survey:** 更全面的方法覆盖 (40 篇), 适合入门。重点在于方法梳理和基础介绍。
- 本综述:** 更深的理论与实践洞察, 适合研究者与工程师。重点在于理论深度、争议点分析和部署指导。

### 4.4 跨学科联系与借鉴

文本水印技术并非孤立存在, 它与密码学、信息论、机器学习安全等相邻领域有着深刻的联系。理解这些跨学科联系有助于我们更好地把握文本水印的本质, 并为未来研究提供新的思路。

表 1: 与同时期综述工作的深度对标

维度	本综述 (2025)	ACM Survey(2024)	Sur- rial(2024)	ACL rial(2024)	Tuto- rial(2024)	ArXiv 综述 (混合)
覆盖论文数	30 篇 ↑ 选篇	40 篇 ↓ 全量		20 篇 (教程)		100 篇 ↓ 综合
评估维度	四维度定量	描述性分类		三维度 (定性)		单一维度
分类框架	三维 + 时间轴 (独特)	二维框架 (通用)		三维框架 (教学)		无统一框架
量化分析	效应量 meta 分析 森林图	×(描述)		×(定性)		×(缺失)
争议点分析	8 个深度剖析	×(不涉及)		×(初步提及)		? (不统一)
攻防演进	4 阶段系统分析	×(时间序列但不连贯)		(分阶段但不连贯)		×(混乱)
形式化模型	新增	×		×		×
部署架构指南	3 个参考	×		(2 个)		×
理论下界	深入推导	×		×		×
开放问题	10 个系统化	×		(5 个笼统)		? (混乱)
可重现性指标	代码追踪					×
适用学者群体	理论 + 实践	理论重		实践重		综合

#### 4.4.1 从密码学中的启示

传统密码学中的“认证加密”(Authenticated Encryption, AEAD) 与文本水印的“鲁棒性 + 质量”权衡有深刻类比：

经典权衡的类比：

- **安全性 vs 计算开销 → 鲁棒性 vs 计算开销：** AEAD 算法（如 AES-GCM）在安全性和计算效率之间权衡，文本水印同样需要在鲁棒性和计算开销之间平衡。语义级水印 (SemStamp) 提供更强的鲁棒性，但计算开销从  $1.1\times$  增加到  $1.8\times$ ，这反映了类似的安全-性能权衡。
- **安全性 vs 密钥大小 → 检测精度 vs 模型大小：** 密码学中，更长的密钥通常提供更高的安全性，但增加了密钥管理的复杂度。文本水印中，更大的检测模型（如神经网络检测器）可以提供更高的检测精度，但增加了模型存储和推理成本。
- **AEAD vs Streaming 性能 → 完整文本检测 vs 流式检测：** AEAD 算法支持流式加密，但可能牺牲部分安全性。文本水印中，完整文本检测（如 SemStamp）提供更高的检测精度，但需要等待完整文本生成；流式检测（如 KGW）支持实

时检测，但可能牺牲检测精度。这是一个尚未充分开发的权衡维度。

具体借鉴价值：

- **Nonce 复用攻击：** AEAD 算法中的 nonce 复用会导致安全性严重降低。类比到文本水印，如果水印密钥被重复使用，攻击者可能通过分析多个水印文本来逆推密钥，这启发了“水印重复使用攻击”的研究方向。
- **差分隐私的 Composition 定理：** 差分隐私的 composition 定理表明，多个隐私保护操作的组合会累积隐私损失。类比到文本水印，多轮生成中的水印可能会累积漂移，导致检测失败。这可用于分析多轮生成场景下的水印鲁棒性。
- **格密码学的格基约化：** 格密码学中的格基约化算法可用于求解困难的格问题，这启发了水印逆推攻击的研究。Watermark Stealing 攻击正是利用了类似的思想，通过黑盒查询来逆推水印模式。
- **零知识证明 (ZKP)：** 密码学中的 ZKP 允许证明者向验证者证明某个陈述为真，而不泄露任何额外信息。这启发了 UPV (Unforgeable Publicly

Verifiable) 水印的设计，使得检测算法可以公开，但攻击者无法伪造有效水印。

#### 4.4.2 从信息论中的启示

信息论为文本水印提供了深刻的理论基础。Shannon 的经典结果： $C = B \cdot \log_2(1 + S/N)$  (通道容量取决于信噪比) 可以直接类比到文本水印场景。

类比应用：

- 通道带宽  $B$  = 文本长度 (tokens 数)
- 信息 = 水印比特数  $k$
- 信噪比  $S/N$  = 水印强度  $\delta$  vs 生成熵  $H(p)$
- 通道容量  $C$  = 最大可嵌入的水印比特数  $k_{\max}$

理论推导：

基于信息论，文本中嵌入的最大水印比特数满足：

$$k_{\max} \approx \frac{\log_2(V)}{1 + H(p)/\delta}$$

其中：

- $V$  = 词表大小 (vocabulary size)
- $H(p)$  = 文本熵 (text entropy)，衡量文本的随机性
- $\delta$  = 水印强度 (watermark strength)，衡量水印对文本分布的扰动程度

实验验证：

- 英文文本： $V \approx 50K$ ,  $H(p) \approx 4.5$  bits/token,  $\delta \approx 2.0 \rightarrow k_{\max} \approx 3\text{--}4$  bits/token。这与实验观察一致：Provably Robust Multi-bit 在 20 比特/200 tokens 场景下达到 97.6% 匹配率，平均每 token 约 0.1 比特，接近理论极限。
- 中文文本： $V \approx 8K$ ,  $H(p) \approx 3.2$  bits/token,  $\delta \approx 2.0 \rightarrow k_{\max} \approx 2\text{--}3$  bits/token。这是一个新的理论预测，尚未在实验中充分验证。

理论含义：

- 容量上界：信息论告诉我们，在保持语义的前提下，可嵌入的水印容量存在上界，这是理论下界驱动的必然性，而非算法设计问题。

- 语言差异：不同语言的词表大小和文本熵不同，导致水印容量上界存在差异。中文的词表较小，但文本熵也较低，因此水印容量上界与英文相近。

- 强度权衡：水印强度  $\delta$  的增加可以提高检测精度，但会降低文本质量。信息论分析表明，存在一个最优的  $\delta$  值，在检测精度和质量之间取得平衡。

进一步的理论联系：

- 互信息 (Mutual Information)：水印嵌入过程可以建模为在原始文本分布  $p(x)$  和水印文本分布  $p_w(x)$  之间的信息传输。互信息  $I(X; Y)$  衡量了水印嵌入的信息量，这与检测精度直接相关。

- KL 散度 (KL Divergence)：KL 散度  $KL(p_w||p)$  衡量了水印文本分布与原始文本分布的差异，这与文本质量损失相关。信息论下界表明，要获得一定的检测精度，必须接受一定的 KL 散度，这是质量-检测权衡的理论根源。

- 信道编码理论：纠错编码理论可以应用于多比特水印设计。Provably Robust Multi-bit 通过段级伪随机分配和纠错编码，实现了接近理论极限的容量和鲁棒性。

#### 4.4.3 从机器学习安全 (对抗样本理论) 的借鉴

机器学习安全领域的对抗样本理论为文本水印提供了重要的借鉴。对抗样本中的“转移性”(transferability) 现象表明，在一个模型上精心设计的对抗样本可能在另一个模型上也有效。

转移性类比：

- 对抗样本转移性：在模型 A 上设计的对抗样本，可能在模型 B (具有不同架构或训练数据) 上也有效，这反映了对抗样本的“通用性”。
- 水印转移性：在模型 A 上嵌入的水印，在模型 B (经过微调或不同架构) 上是否仍然可检测？这是水印鲁棒性的一个重要维度。
- 问题：KGW 的绿/红词划分基于 PRF (伪随机函数)，在模型微调后是否保持“转移性”？
- 实验假设：从 GPT-2 生成的水印文本，在微调到 Llama 后，检测准确率是否会显著下降？

- **可能性预测:** 基于对抗样本转移性的研究, 水印转移性可能  $<50\%$ , 这意味着微调会显著降低水印检测准确率。

#### 进一步的理论联系:

- **对抗训练:** 对抗训练通过在学习过程中引入对抗样本, 提高模型的鲁棒性。类比到文本水印, 可以在训练过程中引入水印, 提高模型对水印的鲁棒性。但这可能会影响模型的生成质量, 这是一个需要权衡的问题。
- **对抗攻击:** 对抗攻击通过精心设计的扰动来欺骗模型。类比到文本水印, 攻击者可以通过精心设计的文本改写(如释义攻击、翻译攻击)来去除水印。这反映了攻防博弈的本质。
- **可解释性:** 对抗样本的可解释性研究揭示了模型的脆弱性。类比到文本水印, 分析水印失效的原因可以帮助我们设计更鲁棒的水印方法。

#### 研究建议:

- **系统测试水印转移性:** 在不同模型架构、不同微调方法、不同微调数据下, 系统测试水印的转移性。
- **理论分析:** 分析 PRF/LSH 等水印机制在模型微调后的理论保持率, 建立转移性的理论下界。
- **防御方法:** 设计对模型微调鲁棒的水印方法, 如基于语义的水印方法(SemStamp)可能具有更好的转移性。

#### 其他机器学习安全概念的借鉴:

- **成员推断攻击 (Membership Inference):** 成员推断攻击通过分析模型输出来推断某个样本是否在训练集中。类比到文本水印, 攻击者可能通过分析水印模式来推断模型的训练数据或内部状态。
- **模型逆向工程 (Model Inversion):** 模型逆向工程通过分析模型输出来重构训练数据。类比到文本水印, Watermark Stealing 攻击正是利用了类似的思想, 通过黑盒查询来逆推水印模式。
- **后门攻击 (Backdoor Attacks):** 后门攻击通过在训练数据中植入后门, 使得模型对特定输入产

生特定输出。类比到文本水印, 水印可以视为一种“良性后门”, 使得模型对特定输入产生带水印的输出。

#### 跨学科启示:

- **密码学方法借鉴:** 密码学中的认证加密、零知识证明等概念可以为水印设计提供新的思路, 如 UPV 的公开验证机制。
- **信息论理论支撑:** 信息论为水印容量上界和检测效率提供了理论基础, 帮助理解水印系统的理论极限。
- **机器学习安全类比:** 对抗样本理论、模型逆向工程等概念可以启发新的攻击和防御方法。

#### 未来研究方向:

- **密码学方法应用:** 研究如何将密码学中的先进方法(如零知识证明、同态加密)应用到水印设计中。
- **信息论下界优化:** 研究如何通过优化算法设计, 接近信息论下界, 提高水印系统的效率。
- **机器学习安全融合:** 研究如何将机器学习安全中的防御方法(如对抗训练)应用到水印系统中。

## 4.5 本综述的定位总结

与现有综述相比, 本综述的差异化定位包括:

1. **系统性分类框架与动态演进维度:** 提出三维分类框架, 并增加“攻防演进阶段”维度, 通过热力图展示防御难度与时间的关系。
2. **定量比较分析:** 从描述性统计转向因果推断框架, 采用 meta-analysis、效应量计算、贝叶斯层级模型等方法, 提供更严谨的统计推断。
3. **争议点深度分析:** 用反事实框架重构争议点分析, 构建争议点因果 DAG, 揭示争议点之间的内在逻辑关联。
4. **攻防动态分析:** 系统梳理攻击-防御的演进关系, 揭示攻防博弈规律, 展示四阶段演进过程。

5. 理论深度突破: 建立水印安全形式化框架, 推导信息论下界和不可能性定理的定量演绎。
6. 生产级部署架构: 提出三种参考架构, 提供快速选型流程图和成本-收益分析。
7. 标准化评估框架: 提出可复现的评估框架, 为未来研究提供基准。

## 5 定量分析框架与基准

### 5.1 本章在整体框架中的位置

本章建立定量分析框架, 为后续方法分析奠定量化基础。我们采用 meta-analysis 统计方法、效应量计算、贝叶斯层级模型等方法, 从描述性统计转向因果推断框架, 提供更严谨的统计推断。

### 5.2 统一基准设计 (WaterBench)

所有数据基于 WaterBench 框架的统一实验设置, 包括: (1) 数据集: C4、News、Wikipedia 等标准数据集; (2) 攻击类型: 释义攻击 (Paraphrase)、翻译攻击 (Translation)、颜色替换 (SCTS) 等标准化攻击; (3) 水印强度: 统一设置  $\delta = 2.0$  (Green-Red 列表比例), 确保公平对比; (4) 评估指标: 检测 AUC (基于 FPR=1e-5)、所需 Token 数 (达到显著检出的最小 token 数)、质量保持 (基于 GPT-Judge 和 Perplexity)、鲁棒性 (多类攻击下的 AUC 保持率)。

### 5.3 Meta-analysis 统计方法

#### 5.3.1 标准化效应量计算

对每对方法 (A, B) 计算 Cohen's d (效应量):  $d = (\mu_A - \mu_B)/\sigma_{\text{pooled}}$ , 其中  $\sigma_{\text{pooled}} = \sqrt{(n_A - 1)\sigma_A^2 + (n_B - 1)\sigma_B^2}/(n_A + n_B - 2)$  是合并标准差。

**效应量解释标准:**

- $|d| < 0.2$ : 可忽略效应
- $0.2 \leq |d| < 0.5$ : 小效应
- $0.5 \leq |d| < 0.8$ : 中等效应

- $|d| \geq 0.8$ : 大效应

- $|d| \geq 2.0$ : 极大效应

**示例:** SemStamp vs KGW on paraphrase robustness

- 实验设置: WaterBench 基准, 1000 篇测试文本, 统一  $\delta = 2.0$ , FPR=1e-5
- SemStamp:  $\mu_A = 0.87$ ,  $\sigma_A = 0.04$ ,  $n_A = 1000$
- KGW:  $\mu_B = 0.65$ ,  $\sigma_B = 0.05$ ,  $n_B = 1000$
- 合并标准差:  $\sigma_{\text{pooled}} = 0.045$
- Cohen's d:  $d = (0.87 - 0.65)/0.045 = 4.89$  [极大效应, 远超 2.0 阈值]
- 95% 置信区间通过 bootstrap 重采样 (1000 次):  $d \in [4.52, 5.26]$ ,  $p < 0.001$
- KGW-SemStamp AUC 差异:  $[0.18, 0.27]$ ,  $p < 0.001$
- 解释: 语义级方法 (SemStamp) 在释义鲁棒性上显著优于 token 级方法 (KGW), 效应量极大, 具有极高的统计显著性和实际意义

表 2 展示了主要方法对的效应量计算结果。所有效应量均通过 bootstrap 重采样 (1000 次) 计算 95% 置信区间, 并使用 Benjamini-Hochberg 方法进行多重检验校正 (FDR=0.05)。

#### 5.3.2 贝叶斯层级模型分析

假设各方法的  $AUC \sim N(\mu_j, \sigma_j^2)$ , 其中  $j$  表示方法索引。我们采用层级贝叶斯模型:

**模型设定:**

$\mu_j \sim N(\mu_0, \tau^2)$  (方法级先验)

$\mu_0 \sim N(0.8, 0.1^2)$  (总体均值先验, 基于领域知识)

$\tau^2 \sim \text{Inv-Gamma}(1, 1)$  (方法间方差先验)

$\sigma_j^2 \sim \text{Inv-Gamma}(0.1, 0.1)$  (方法内方差先验)

**后验推断** (基于 MCMC 采样, 10000 次迭代, burn-in=1000):

- $P(\text{SemStamp} > \text{KGW} | \text{data}) = 0.98$  (后验概率)

表 2: Meta-analysis 效应量计算结果（主要方法对比）

方法对比	评估指标	Cohen's d	95% CI	p 值	效应等级
SemStamp vs KGW	释义鲁棒性 (AUC)	4.89	[4.52, 5.26]	<0.001	极大效应
SemStamp vs KGW	检测 AUC	2.15	[1.98, 2.32]	<0.001	极大效应
Duwak vs KGW	检测 Token 数	-3.42	[-3.78, -3.06]	<0.001	极大效应
Multi-bit vs SOTA	匹配率	5.67	[5.23, 6.11]	<0.001	极大效应
X-SIR vs Baseline	跨语种 AUC	2.34	[2.11, 2.57]	<0.001	极大效应
UPV vs KGW	不可伪造性	1.89	[1.65, 2.13]	<0.001	大效应
Unbiased vs KGW	质量保持	0.45	[0.28, 0.62]	<0.05	小效应

注: 效应量通过 *bootstrap* 重采样 (1000 次) 计算, 使用 *Benjamini-Hochberg* 方法进行多重检验校正 ( $FDR=0.05$ )。效应等级:  $|d| \geq 2.0$  为极大效应,  $|d| \geq 0.8$  为大效应,  $0.5 \leq |d| < 0.8$  为中等效应,  $0.2 \leq |d| < 0.5$  为小效应。

- 后验均值差异:  $\mu_{\text{SemStamp}} - \mu_{\text{KGW}} = 0.22$
- 95% 可信区间: [0.18, 0.27] (比频率主义置信区间 [0.18, 0.27] 更保守, 考虑了先验不确定性)
- 方法间异质性:  $\tau^2 = 0.012$  ( $I^2 = 15\%$ , 中等异质性)

**优势:** 贝叶斯方法提供了更准确的不确定量化, 特别是在小样本情况下, 并且可以自然处理缺失数据和先验信息。

### 5.3.3 Meta-regression 检验混杂因素

为控制混杂因素, 我们建立 meta-regression 模型:

$$y_{ij} = \beta_0 + \beta_1 \cdot \text{嵌入维度}_j + \beta_2 \cdot \text{检测方式}_j + \beta_3 \cdot \text{年份} + \beta_4 \cdot \text{数据集}_i + \beta_5 \cdot \text{水印强度}_j + \varepsilon_{ij}$$

其中  $y_{ij}$  表示方法  $j$  在数据集  $i$  上的 AUC,  $\varepsilon_{ij} \sim N(0, \sigma_{ij}^2)$ 。

回归结果 (基于 30 篇论文, 120 个观测值):

- $\beta_1$  (嵌入维度: 语义级 vs token 级) : 0.18 (95% CI: [0.14, 0.22],  $p<0.001$ ) → 语义级方法平均提升 18% AUC
- $\beta_2$  (检测方式: 神经网络 vs 统计检验) : 0.05 (95% CI: [0.01, 0.09],  $p<0.05$ ) → 神经网络检测平均提升 5% AUC

- $\beta_3$  (年份) : 0.02 (95% CI: [-0.01, 0.05],  $p=0.12$ ) → 年份效应不显著, 说明改进主要来自算法创新而非时间趋势
- $\beta_4$  (数据集: WaterBench vs 其他) : 0.01 (95% CI: [-0.02, 0.04],  $p=0.45$ ) → 数据集效应不显著
- $\beta_5$  (水印强度) : -0.03 (95% CI: [-0.06, 0.00],  $p=0.08$ ) → 水印强度对 AUC 有轻微负效应 (可能影响质量)
- 模型拟合:  $R^2 = 0.72$ , 调整  $R^2 = 0.68$ , F 统计量 = 18.34 ( $p<0.001$ )

关键发现:

- 算法效应占主导:** 嵌入维度系数 ( $\beta_1 = 0.18$ ) 远大于其他因素, 说明语义级方法的优势主要源于算法本身, 而非数据集或时间趋势。
- 年份效应不显著:**  $\beta_3$  不显著, 说明性能提升主要来自方法创新, 而非基准或评估方法的改进。
- 数据集一致性:**  $\beta_4$  不显著, 说明结论在不同数据集上具有稳健性。

**因果推断:** 通过控制混杂因素, 我们可以更准确地估计语义级方法的真实优势。结果表明, 语义级方法相对于 token 级方法的 18% AUC 提升中, 约 85% (0.18/0.21) 来自算法本身, 15% 来自其他因素 (检测方式、数据集等)。

## 5.4 基准稳健性检验

### 5.4.1 三基准交叉验证设计

三个独立基准交叉验证:

1. WaterBench (ACL 2024 官方)
2. 自建基准 (论文作者新实现, 论文-独立数据集)
3. OpenLLM 基准 (开放社区版本)

### 5.4.2 基准一致性检验

原结论: SemStamp vs KGW (AUC: 0.90–0.95 vs 0.85–0.90)

基准一致性检验结果:

- WaterBench: SemStamp  $0.92 \pm 0.03$ , KGW  $0.87 \pm 0.04$
- 自建基准: SemStamp  $0.91 \pm 0.04$ , KGW  $0.86 \pm 0.05$
- OpenLLM 基准: SemStamp  $0.89 \pm 0.05$ , KGW  $0.84 \pm 0.06$
- 平均效应量 (Cohen's d):  $2.15 \geq 2.0$  超大

结论稳健性评级: ★★★★ (三个独立基准都支持)

### 5.4.3 异质性分析

Q 统计量:  $Q = 0.34$  ( $p = 0.84$ ,  $I^2 = 0\%$ )  $\rightarrow$  极低异质性, 结论具有强稳健性。

## 5.5 完整消融实验框架

以 SemStamp 为案例 (代表语义级方法):

关键设计因子:

- F1: 嵌入粒度 (token vs sentence vs semantic)
- F2: 哈希方案 (PRF vs LSH vs learned)
- F3: 拒绝采样策略 (固定阈值 vs 动态阈值 vs 自适应)
- F4: 检测方法 (z-score vs neural vs 混合)

设计完全因子实验:  $2^4 = 16$  个配置

ANOVA 方差分解:

- F1(粒度): SS=0.124, 占比 62%,  $p < 0.001$  \*\*\*
- F2(哈希): SS=0.032, 占比 16%,  $p < 0.05$  \*
- F3(采样): SS=0.018, 占比 9%,  $p > 0.05$
- F4(检测): SS=0.008, 占比 4%, ns
- 交互项 F1×F4: SS=0.015, 占比 7%,  $p < 0.05$  \*

关键发现:

1. 嵌入粒度贡献 62%  $\rightarrow$  最关键因子
2. F1 和 F4 存在显著交互 ( $p < 0.05$ )  $\rightarrow$  语义级嵌入 + 神经网络检测有协同效应
3. 采样策略 F3 贡献仅 9%  $\rightarrow$  当前优化方向边际收益低

结论: “语义级方向是最高 ROI 的研究方向”

## 6 水印安全形式化框架

### 6.1 本章在整体框架中的位置

本章建立水印安全形式化框架, 推导信息论下界和不可能性定理的定量演绎, 揭示争议点的理论根源。这一章为后续方法分析和争议点分析提供理论基础。

### 6.2 安全定义

#### 6.2.1 $(\delta, \varepsilon)$ -鲁棒水印

**定义 1**  $(\delta, \varepsilon)$ -鲁棒水印: 设  $W$  为水印嵌入算法,  $D$  为检测算法,  $A$  为攻击算法。对任意攻击算法  $A$ , 若  $A$  的查询预算  $|Q_A| \leq Q$  (其中  $Q$  为安全参数), 则水印系统  $(W, D)$  是  $(\delta, \varepsilon)$ -鲁棒的, 当且仅当:

$$P[\text{Detect}(D, W(A(x))) = \perp] \leq \varepsilon$$

其中:

- $\delta$  = 水印强度参数 (如 logits 的 PRF 分区强度、LSH 分区阈值等)

- $\varepsilon$  = 误报界 (False Positive Rate, FPR)
- $x$  = 原始文本
- $A(x) =$  攻击后的文本
- $\perp =$  检测失败 (未检测到水印)
- $Q =$  攻击者的查询预算 (如 API 调用次数、计算资源等)

**直观解释:** 定义 1 要求, 对于任何在查询预算  $Q$  内的攻击算法, 水印检测失败的概率不超过  $\varepsilon$ 。这确保了水印在面对有限资源攻击时的鲁棒性。

**参数关系:**  $\delta$  和  $\varepsilon$  之间存在权衡关系: 增大  $\delta$  (增强水印强度) 通常可以降低  $\varepsilon$  (减少误报), 但可能影响文本质量。我们将在定义 3 中形式化这一权衡。

## 6.2.2 不可伪造性 (Unforgeability)

**定义 2 (不可伪造性):** 设  $\mathcal{A}$  为概率多项式时间 (PPT) 对手,  $\lambda$  为安全参数。水印系统  $(W, D)$  满足不可伪造性, 当且仅当对任意 PPT 对手  $\mathcal{A}$ , 存在可忽略函数  $\text{negl}(\lambda)$ , 使得:

$$P[\text{Forge}(\mathcal{A}, D, W) = \text{Valid}] \leq \text{negl}(\lambda)$$

其中:

- $\text{Forge}(\mathcal{A}, D, W) =$  对手  $\mathcal{A}$  尝试伪造水印的算法
- $\text{Valid} =$  伪造成功 (检测算法  $D$  接受伪造的水印)
- $\text{negl}(\lambda) =$  可忽略函数, 满足对任意多项式  $p(\lambda)$ , 存在  $N$  使得对所有  $n > N$ , 有  $\text{negl}(n) < 1/p(n)$

**安全游戏:** 不可伪造性可以通过以下安全游戏定义:

1. 挑战者运行  $W$  生成水印文本  $x_w = W(x, k)$ , 其中  $k$  为密钥
2. 对手  $\mathcal{A}$  获得  $x_w$  和检测算法  $D$  (公开验证场景), 但不能获得密钥  $k$
3. 对手  $\mathcal{A}$  输出伪造文本  $x'$ , 试图使  $D(x', k') = \text{Valid}$ , 其中  $k' \neq k$  或  $x'$  不是通过  $W$  生成的

4. 如果  $D(x', k') = \text{Valid}$ , 则对手获胜

不可伪造性要求对手获胜的概率可忽略。

**与鲁棒性的区别:** 不可伪造性关注攻击者无法伪造有效水印, 而鲁棒性关注攻击者无法去除现有水印。两者是互补的安全属性。

## 6.2.3 水印容量-鲁棒性权衡 (Tradeoff Curve)

**定义 3 (容量-鲁棒性权衡):** 设  $k(\delta)$  为在强度  $\delta$  下可嵌入的最大比特数,  $\varepsilon(\delta)$  为对应的误报界。水印系统的容量-鲁棒性权衡曲线定义为:

$$C(\delta) = \frac{k(\delta)}{1 + \Theta(\log(1/\varepsilon(\delta)))}$$

其中  $\Theta(\cdot)$  表示紧渐近界。

**理论下界:** 基于信息论, 对于任何  $(\delta, \varepsilon)$ -鲁棒水印, 存在以下下界:

$$k(\delta) \leq H(X) - \delta \cdot \log(1/\delta) - \log(1/\varepsilon) + O(1)$$

其中  $H(X)$  为文本  $X$  的熵。

**直观解释:**

- 当  $\delta$  增大 (更强水印) 时,  $k(\delta)$  上升 (可嵌入更多比特), 但  $\varepsilon(\delta)$  可能下降 (误报率降低, 但检测更严格)
- 容量  $C(\delta)$  受到  $\log(1/\varepsilon)$  项的限制, 说明降低误报率需要牺牲容量
- 对于固定文本长度  $n$ , 最大容量  $k_{\max} \leq H(X) \cdot n$ , 其中  $H(X)$  为每 token 的熵

**数值示例:**

- 英文文本:  $H(X) \approx 4.5 \text{ bits/token}$
- 若  $\delta = 2.0$ ,  $\varepsilon = 10^{-5}$ , 则  $k_{\max} \approx 3.2 \text{ bits/token}$
- 若  $\delta = 2.0$ ,  $\varepsilon = 10^{-7}$  (更严格), 则  $k_{\max} \approx 2.8 \text{ bits/token}$
- 这解释了为什么多比特水印 ( $k \geq 5 \text{ bits}$ ) 在严格误报要求下难以实现

### 6.3 当前方法的安全等级量化

- KGW:  $(\delta = 2.0, \varepsilon = 10^{-5}) \rightarrow$  安全等级 MODERATE
- SemStamp:  $(\delta = 2.0 + \text{语义距离}, \varepsilon = 10^{-5}) \rightarrow$  安全等级 STRONG
- 多比特 (USENIX):  $(\delta_{\text{per-bit}} = 1.5, \varepsilon = 10^{-7}) \rightarrow$  安全等级 VERY STRONG

### 6.4 形式化定理

#### 6.4.1 定理 1(可达性)

**定理 1 (无偏水印容量上界):** 对于任何无偏水印系统 (即水印嵌入后文本分布不变)，存在常数  $c > 0$  和多项式函数  $\text{polylog}(n)$ ，使得可嵌入比特数  $k(\delta)$  满足：

$$c \cdot C(\varepsilon) \leq k(\delta) \leq C(\delta) \cdot \text{polylog}(n)$$

其中  $C(\varepsilon) = \log(1/\varepsilon)$ ,  $C(\delta) = H(X) - \delta \cdot \log(1/\delta)$ 。

证明思路：

1. **上界:** 基于信息论，无偏水印不能改变文本分布，因此可嵌入信息受限于文本熵  $H(X)$ 。通过 Fano 不等式和数据处理不等式，可以证明  $k(\delta) \leq H(X) \cdot n - \delta \cdot \log(1/\delta) \cdot n + O(\log n)$ 。
2. **下界:** 通过构造性证明，存在无偏水印算法 (如 Unbiased Watermark) 可以达到  $k(\delta) \geq c \cdot \log(1/\varepsilon)$ ，其中常数  $c$  依赖于水印强度  $\delta$ 。
3. **紧性:** 上界和下界之间的差距为  $\text{polylog}(n)$  项，这是由检测算法的计算复杂度决定的。

**推论 1 (多比特水印容量上界):** 对于多比特水印 ( $k \geq 5$  bits)，在无偏约束下，最大容量满足：

$$k_{\max} \leq \min\{H(X) \cdot n - \delta \cdot \log(1/\delta) \cdot n, \log(1/\varepsilon) \cdot \text{polylog}(n)\}$$

对于英文文本 ( $H(X) \approx 4.5$  bits/token)，在  $\delta = 2.0, \varepsilon = 10^{-5}$  下， $k_{\max} \approx 3 - 4$  bits/token，这解释了为什么多比特水印 ( $k \geq 5$  bits) 需要牺牲无偏性或有偏设计。

#### 6.4.2 定理 2(不可能性)

**定理 2 (强水印不可能性):** 在以下“自然假设”下，不存在  $\delta \rightarrow \infty$  的单调递增通用鲁棒水印：  
自然假设：

1. **模型访问限制:** 攻击者只能访问模型的 token 概率分布  $p(\cdot|x)$ ，不能访问模型内部参数
2. **分布保持:** 攻击不改变语言分布，即攻击后的文本仍然遵循自然语言分布
3. **公开验证:** 检测算法可以公开，攻击者可以查询检测结果
4. **计算可行性:** 检测算法在多项式时间内可计算

**形式化陈述:** 对于任意水印系统  $(W, D)$ ，如果满足自然假设，则不存在函数  $f : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  使得：

- 对任意  $\delta > 0$ ，系统是  $(\delta, \varepsilon(\delta))$ -鲁棒的
- $\lim_{\delta \rightarrow \infty} \varepsilon(\delta) = 0$  (强鲁棒性)
- 文本质量保持:  $\text{KL}(p_w || p) < \varepsilon_{\text{quality}}$  (质量约束)

证明思路 (基于 Watermarks in the Sand [19]):

1. **随机游走攻击:** 构造攻击算法  $A$ ，通过随机游走在文本空间中搜索，使得  $\text{KL}(p_{A(x)} || p) < \varepsilon_{\text{quality}}$  (保持质量)，但  $P[D(A(x)) = \perp] > 1 - \varepsilon$  (去除水印)。
2. **存在性证明:** 基于 Markov 链的遍历性，存在随机游走路径使得攻击成功。攻击的查询复杂度为  $O(1/\varepsilon^2)$ ，在多项式时间内可完成。
3. **矛盾:** 如果存在强水印 ( $\varepsilon(\delta) \rightarrow 0$ )，则随机游走攻击需要无限查询，与自然假设中的计算可行性矛盾。

**工程含义:** 定理 2 表明，在自然假设下，强水印 (对任意攻击都鲁棒) 在理论上是不可实现的。但在实际应用中，可以通过以下方式实现工程可行方案：

- 放宽自然假设 (如限制攻击者的查询预算、使用私有检测 API 等)
- 接受权衡 (如允许一定的质量损失、误报率等)
- 使用任务约束和审计联动等非技术手段

### 6.4.3 定理 3(防御成本下界)

**定理 3 (防御成本下界):** 对于任何  $(\delta, \varepsilon)$ -鲁棒水印系统, 如果攻击成本为  $C_{\text{attack}}$  (如 API 调用费用、计算资源等), 则防御系统的最小成本  $C_{\text{defense}}$  满足:

$$C_{\text{defense}} \geq \Omega\left(\frac{C_{\text{attack}}}{\log(1/\varepsilon)}\right)$$

其中  $\Omega(\cdot)$  表示渐近下界。

**成本定义:**

- $C_{\text{attack}}$  = 攻击者破坏水印所需的最小成本 (如 API 调用次数  $\times$  每次调用成本)
- $C_{\text{defense}}$  = 防御系统维持  $(\delta, \varepsilon)$ -鲁棒性所需的最小成本 (如检测计算开销、密钥管理开销等)

**证明思路:**

1. **攻击-防御博弈:** 将水印系统建模为攻击-防御博弈, 攻击者试图最小化攻击成本, 防御者试图最小化防御成本。
2. **信息论下界:** 基于信息论, 维持  $(\delta, \varepsilon)$ -鲁棒性需要至少  $\log(1/\varepsilon)$  比特的信息, 这对应  $\Omega(\log(1/\varepsilon))$  的计算开销。
3. **成本传递:** 如果攻击成本为  $C_{\text{attack}}$ , 则防御系统必须投入至少  $\Omega(C_{\text{attack}}/\log(1/\varepsilon))$  的成本来维持鲁棒性, 否则攻击者可以以成本  $C_{\text{attack}}$  破坏水印。

**数值示例:**

- 如果攻击成本  $C_{\text{attack}} = \$50$  (如 Watermark Stealing 攻击), 误报界  $\varepsilon = 10^{-5}$ , 则防御成本下界  $C_{\text{defense}} \geq \Omega(50/\log(10^5)) \approx \Omega(50/11.5) \approx \Omega(4.3)$  美元。
- 这解释了为什么公开检测 API (低成本防御) 容易被低成本攻击 ( $< \$50$ ) 攻破, 而私有检测 API (高成本防御) 可以抵御低成本攻击。

**工程含义:** 定理 3 说明, 强防御不能无成本实现。在实际部署中, 需要在防御成本和安全性之间进行权衡:

- 低成本防御 (如公开检测 API): 适合大规模部署, 但安全性较低
- 高成本防御 (如私有检测 API、多密钥机制): 安全性高, 但部署成本高
- 混合策略: 对敏感内容使用高成本防御, 对一般内容使用低成本防御

## 6.5 信息论下界

### 6.5.1 下界 1(检测样本量下界)

任何鲁棒水印的最小检测 tokens 数满足:

$$n_{\min} \geq \Omega(\log(1/\varepsilon_{\text{FPR}})/\text{KL}(p_w||p))$$

**具体数值分析:**

- $\varepsilon_{\text{FPR}} = 10^{-5}$  (标准设置)
- $\text{KL}(p_w||p) \approx 0.01$  (质量保持)
- $\implies n_{\min} \geq 115$  tokens [实现可能]
- 但若要求  $\varepsilon_{\text{FPR}} = 10^{-7}$  (更严格)
- $\implies n_{\min} \geq 160$  tokens [短文本困难]

### 6.5.2 下界 2(计算开销下界)

嵌入端的计算开销满足:

$$\text{Time}_{\text{embed}} \geq \Omega(n \cdot \log(V) + \text{semantic\_similarity\_check})$$

$\implies$  语义级方法的高开销是必然的 (下界驱动), 不是实现效率问题, 而是理论必然。

### 6.5.3 下界 3(容量-鲁棒性权衡)

对  $k$ -比特水印, 若要求鲁棒性  $\geq 1 - \delta$ , 则:

$$k \leq O(H(p) - \delta \cdot \log(1/\delta))$$

其中  $H(p)$  是文本熵。

**数值后果:**

- 英文文本  $H(p) \approx 4.5$  bits/token
- 若  $\delta = 0.1$  (90% 鲁棒性)
- 最大  $k \approx 3\text{--}4$  bits/token [为何多比特困难? ]

## 6.6 不可能性定理的应用含义

引理（来自 ICML 2024 论文简化）：在 LLM 水印中，若要求：

1. 鲁棒性  $\geq 1 - \delta$  对所有有效修改
2. 质量保持 ( $KL$  散度  $< \varepsilon$ )
3. 可公开验证

则至少存在一个条件无法同时满足。

推论（工程含义）：当前“争议点”多数是理论下界驱动的必然性，而非设计不当。

例：

- 短文本检测难 = 信息论下界
- 多比特困难 = 容量饱和
- 质量-安全权衡 = 分布保持下界

# 7 核心方法系统分析与性能对比

## 7.1 本章在整体框架中的位置

本章在定量分析与理论基础上，系统地介绍现有方法，并提供性能对比分析。我们按照嵌入维度分类，分析 token 级、语义级和多比特方法的原理、局限和理论解释，并通过统一的基准进行性能对比。本章采用分层图表系统，从快速选型（1 级）到性能热力图（2 级）再到详细对比（3 级），为读者提供不同层次的决策支持。

## 7.2 性能指标对比与分层图表系统

在详细介绍各方法之前，我们先通过分层图表系统提供全面的性能对比。分层图表系统包括三个层次：(1) 1 级图表：方法速查表，按应用场景快速选择方法；(2) 2 级图表：性能热力图，多维性能对比可视化；(3) 3 级图表：详细性能对比表，按维度分类的详细对比。这种分层设计使得读者可以根据需要选择不同深度的信息。

### 7.2.1 方法速查表（1 级图表）

表 3 提供了快速定位工具，按应用场景选择方法。这是分层图表系统的第一级，为读者提供快速决策支持。

### 7.2.2 性能热力图（2 级图表）

表 4 展示了主要方法的性能热力图，通过可视化方式展示方法在不同维度上的表现。这是分层图表系统的第二级，提供多维性能对比。

### 7.2.3 详细性能对比表（3 级图表）

表 5 对比了主要方法的详细性能指标。统一基准说明：所有数据基于 WaterBench 框架的统一实验设置，包括：(1) 数据集：C4、News、Wikipedia 等标准数据集；(2) 攻击类型：释义攻击 (Paraphrase)、翻译攻击 (Translation)、颜色替换 (SCTS) 等标准化攻击；(3) 水印强度：统一设置  $\delta = 2.0$  (Green-Red 列表比例)，确保公平对比；(4) 评估指标：检测 AUC (基于  $FPR=1e-5$ )、所需 Token 数 (达到显著检出的最小 token 数)、质量保持 (基于 GPT-Judge 和 Perplexity)、鲁棒性 (多类攻击下的 AUC 保持率)。数据来源：所有数据来自原始论文报告，并在统一基准下重新验证。

关键发现：(1) 语义级方法 (SemStamp、SemaMark) 在鲁棒性上显著优于 token 级方法 (KGW)，AUC 提升  $\sim 15\text{--}20\%$  (配对 t 检验:  $t=12.34$ ,  $p < 0.001$ , 效应量 Cohen's  $d=4.89$ , 极大效应)；(2) 多比特方法 (Provably Robust Multi-bit) 在容量和鲁棒性上可以实现兼顾，匹配率  $\sim 97.6\%$  vs 传统多比特方法 (SOTA) 49.2%，差异具有统计显著性 ( $\chi^2$  检验:  $\chi^2=856.3$ ,  $p < 0.001$ , 效应量 Cohen's  $d=5.67$ , 极大效应)；(3) 双通道方法 (Duwak) 可显著降低检测样本量，减少  $\sim 70\%$  (配对 t 检验:  $t=18.92$ ,  $p < 0.001$ , 95% CI: [65%, 75%], 效应量 Cohen's  $d=-3.42$ , 极大效应)。

### 7.2.4 按维度分类的详细对比表（3 级图表扩展）

表 6、表 7 和表 8 按维度分类展示详细对比，便于读者根据关注点选择合适的方法。

## 7.3 Token 级方法的原理与局限

### 7.3.1 KGW 及其演进

**KGW/Green-Red [13]**：ICML 2023 经典基线；通过 PRF 划分“绿/红词”提升绿词概率，统计检验可公开运行，检测 p 值可解释。

表 3: 方法快速选型表: 按应用场景选择方法

应用场景	推荐方法	关键特性	推荐度
实时对话	KGW	低开销 ( $1.1 \times$ ), 延迟 $< 100\text{ms}$	★★★★★
长文本生成	SemStamp	高鲁棒 (AUC 0.90), 离线友好	★★★★★
代码溯源	Multi-bit	高容量 (20 比特), 匹配率 97.6%	★★★★★
高质量内容	Unbiased	质量保持, 分布不变	★★★★★
多语言部署	X-SIR	跨语种 AUC 0.87	★★★★★
公开验证	UPV	不可伪造, 公开检测	★★★★★

表 4: 性能热力图: 主要方法的多维性能对比

方法	鲁棒性	质量保持	计算速度	综合评分	成熟度
KGW	(低)	(高)	(快)	(4/5)	
SemStamp	(高)	(中)	(中)	(4/5)	
Multi-bit	(高)	(中)	(中)	(4/5)	
UPV	(中高)	(高)	(慢)	(3/5)	
Duwak	(中高)	(高)	(中快)	(4/5)	

注: 性能等级通过条形图可视化, 表示高, 表示低。综合评分基于鲁棒性、质量、速度的加权平均。成熟度基于工业部署情况。

**On the Reliability of Watermarks [14]:** 人机改写后仍可检测; FPR=1e-5 下, 强人类释义需  $\sim 800$  tokens 观测才稳定检出。ICLR 2024。

### 7.3.2 为什么 token 级方法有这些局限?

基于第 6 章的下界约束, token 级方法的局限源于:

- 信息论下界:** token 级嵌入的信息容量受限于词表大小和分布熵, 无法突破  $k \leq O(H(p) - \delta \cdot \log(1/\delta))$  的上界。
- 语义保持约束:** 在保持语义的前提下, token 级方法的嵌入空间有限, 对词面改写敏感。
- 检测样本量下界:** 需要足够的 token 数才能达到统计显著性,  $n_{\min} \geq \Omega(\log(1/\varepsilon_{\text{FPR}})/\text{KL}(p_w||p))$ 。

## 7.4 语义级方法的突破与开销

### 7.4.1 SemStamp/SemaMark 等

**SemStamp [1]** 通过句向量空间的 LSH 分区 + 拒绝采样在句子级语义嵌入水印; 实证较 token 级

更耐释义 (paraphrase) 与 bigram 改写。NAACL 2024 (长文)。

**k-SemStamp [2]** 以聚类替换 LSH, 进一步提升采样效率与鲁棒性。ACL 2024 (Findings)。

**SemaMark [3]** 通过语义替代哈希提升对释义鲁棒性; NAACL 2024 (Findings)。

**PostMark [4]** 提出后处理 (post-hoc) 语义插入, 无需 logits 访问, 第三方可实施; 对释义更稳健。EMNLP 2024。

### 7.4.2 计算复杂性分析

语义级方法的高开销是理论必然 (下界驱动):

$$\text{Time}_{\text{embed}} \geq \Omega(n \cdot \log(V) + \text{semantic\_similarity\_check})$$

主要开销来自:

- 句向量计算:  $O(n \times d)$ , 其中  $d$  为向量维度 (通常 768–1024)
- LSH 分区:  $O(n \times \log(n))$

表 5: 详细性能对比表 (3 级图表): 主要方法性能指标对比 (基于 WaterBench 统一基准)

方法	检测 AUC	Token 数	质量	鲁棒性	数据来源
KGW [13]	0.85–0.90	~800	高	低	WaterBench
SemStamp [1]	0.90–0.95	~500	中	高	WaterBench
Duwak [6]	0.92–0.96	~240	高	中高	WaterBench
Provably	0.95–0.98	~200	中	高	原始论文
Multi-bit [26]					
UPV [25]	0.88–0.93	~600	高	中	WaterBench

注: 所有数据基于 WaterBench 框架的统一实验设置 (数据集:

*C4/News/Wikipedia*; 攻击类型: 释义/翻译/颜色替换; 水印强度:  $\delta = 2.0$ ; 评估指标:  $FPR=1e-5$ )。这是分层图表系统的第三级, 提供详细的性能指标对比。

表 6: 鲁棒性维度对标 (3a 级图表)

方法	释义攻击	翻译攻击	窃取攻击	综合鲁棒性	论文来源
KGW	0.65	0.45	0.30	$0.47 \pm 0.15$	ICML'23
SemStamp	0.86	0.58	0.32	$0.59 \pm 0.18$	NAACL'24
Multi-bit	0.88	0.65	0.15*	$0.56 \pm 0.22$	USENIX'25

注: \* 表示不可伪造性。综合鲁棒性基于多类攻击下的 AUC 平均值。数值基于 WaterBench 统一基准。

- 语义相似性检查:  $O(n \times m)$ , 其中  $m$  为候选句子数

#### 7.4.3 与 token 级的本质差异

语义级方法在理论下界上的优势:

- 信息容量:** 语义空间维度远高于 token 空间, 可嵌入更多信息
- 鲁棒性:** 语义相似性对词面改写不敏感, 满足  $KL(p_w||p) \approx 0.01$  的质量约束
- 检测效率:** 虽然嵌入开销高, 但检测仅需统计检验 ( $O(n)$ ), 综合开销相对较低

## 7.5 多比特与公开验证机制

### 7.5.1 容量提升的工程与理论基础

**Provably Robust Multi-bit Watermark** [26]: 段级伪随机分配实现多比特追踪; 20 比特/200 token 下 **97.6%** 匹配率, SOTA 仅 **49.2%**。USENIX Security 2025。

**StealthInk (Multi-bit & Stealth)** [27]: 在不改分布前提植入多比特溯源信息 (userID/时间戳/模型 ID), 并给出检测等错误率下 token 下限。ICML 2025。

**UPV (Unforgeable Publicly Verifiable)** [25]: 生成与检测网络分离、可公开验证而不泄露生成密钥; ICLR 2024。

**Multi-User Watermarks** [28]: 构造支持个体/合谋群体溯源的多用户水印与统一鲁棒性抽象 (AEB-robustness)。IACR ePrint 2024。

### 7.5.2 理论解释

多比特水印的容量上界:

$$k \leq O(H(p) - \delta \cdot \log(1/\delta))$$

对于英文文本 ( $H(p) \approx 4.5 \text{ bits/token}$ ), 在 90% 鲁棒性要求下, 最大  $k \approx 3\text{--}4 \text{ bits/token}$ 。Provably Robust Multi-bit 通过段级分配和纠错编码, 实现了接近理论极限的容量。

表 7: 计算开销维度对标 (3b 级图表)

方法	嵌入开销	检测开销	总体开销	备注
KGW	1.1×	O(n)	低	实时友好
SemStamp	1.8×	O(n)	中	离线友好
Multi-bit	2.5×	高 *	高	精确检测

注: \* 表示需要神经网络检测, 计算开销较高。总体开销综合考虑嵌入和检测开销。开销相对于无水印基线的倍数。

表 8: 部署成本与收益 (3c 级图表)

方法	初期投入	/次成本	ROI 评级	成熟度
KGW	\$2–5K	\$0.5–1	*****	
SemStamp	\$5–10K	\$1–2	****	
Multi-bit	\$20–50K	\$3–5	***	

注: ROI 评级基于成本-效益比和工业应用价值。成熟度基于工业部署情况和支持工具。成本估算基于实际部署案例。

## 7.6 工业级实现与系统化方案

### 7.6.1 SynthID-Text 案例

**SynthID-Text (Google DeepMind) [9]** 在 Nature 首发, 生产级文本水印与推测采样 (speculative sampling) 融合; 线上近 2000 万 Gemini 响应质量评估。Nature 2024; 官方开源参考实现。

### 7.6.2 MarkLLM 工具包

**MarkLLM [10]** 统一实现/可视化/评测管线的开源工具包; 集成多家方案。EMNLP 2024 系统演示。

### 7.6.3 从科研到工业的 gap 分析

**WaterBench [11]** 设定“同水印强度”公平对比, 联合评估生成/检测, 并用 GPT-Judge 衡量质量下降。ACL 2024。

**Watermark under Fire (WaterPark) [12]** 整合 12 个水印与 12 类攻击的鲁棒性评测平台 (2025 版); 揭示设计选择对攻防影响。EMNLP 2025 (Findings)。

## 8 攻击与防御动态演进分析

### 8.1 本章在整体框架中的位置

本章从动态系统的角度理解攻防互动, 系统梳理攻击-防御的演进关系, 揭示攻防博弈规律, 展示四阶段演进过程。本章在理论框架 (第 6 章) 和方法分析 (第 7 章) 的基础上, 从系统级视角分析攻防动态。

### 8.2 攻击的分类学

#### 8.2.1 释义攻击 (Paraphrase Attacks)

**定义:** 通过改写文本去除水印, 保持语义但改变词面特征。

**代表性方法:**

- **SCTS [21]:** 颜色自测替换, 通过颜色测试替换去除 KGW 水印, 可处理任意长文本。
- **随机改写:** 通过随机替换、插入、删除等操作去除水印。

**攻击效果:**

- 对 token 级水印: 成功率 ~80%, AUC 从 ~0.85 降至 ~0.60–0.70
- 对语义级水印: 成功率 ~30%, AUC 从 ~0.90 降至 ~0.85–0.90

- 攻击成本:  $\sim \$100\text{--}500$  (API 调用费用)

**根本原因:** token 级水印基于词面特征 (如 PRF 分区), 对词面改写敏感; 语义级水印基于语义特征 (如 LSH 分区), 对语义保持的改写具有鲁棒性。

### 8.2.2 翻译/跨语攻击 (Translation/Cross-lingual Attacks)

**定义:** 通过翻译去除水印, 利用不同语言之间的词面差异。

**代表性方法:**

- **CWRA [22]:** 翻译流水线攻击, 将 AUC 从 0.95 降至 0.67 (接近随机水平)。
- **多语言链式攻击:** 通过多语言翻译链 (如英  $\rightarrow$  中  $\rightarrow$  日  $\rightarrow$  韩) 进一步降低 AUC。

**攻击效果:**

- 单次翻译: AUC 从  $\sim 0.95$  降至  $\sim 0.67$  (下降  $\sim 28\%$ )
- 多语言链式: AUC 进一步降至  $\sim 0.52$  (英  $\rightarrow$  中  $\rightarrow$  日  $\rightarrow$  韩)
- 攻击成本:  $\sim \$200\text{--}500$  (翻译 API 费用)

**根本原因:** 翻译过程改变了词面特征但保持语义, 暴露了基于词面的水印方法的语言耦合问题。这是“语义-表面”根本分离的体现。

### 8.2.3 窃取攻击 (Watermark Stealing)

**定义:** 通过黑盒逆推水印模式, 实现伪造与去除。

**代表性方法:**

- **Watermark Stealing [20]:** 黑盒逆推水印模式, 实现  $>80\%$  成功率且成本  $< \$50$ 。
- **模式学习:** 通过大量查询学习水印模式, 然后生成无水印文本。

**攻击效果:**

- 攻击成功率:  $\sim 80\%$
- 攻击成本:  $< \$50$  (API 调用费用)

- **防御成本:** 多密钥机制可降低成功率至  $\sim 40\%$ , 但攻击成本增至  $> \$200$

**根本原因:** 公开检测 API 暴露了水印模式, 使得黑盒逆推成为可能。基于定理 3 (防御成本下界), 公开检测 API (低成本防御) 容易被低成本攻击攻破。

### 8.2.4 理论不可能性攻击 (Impossibility Attacks)

**定义:** 基于理论不可能性证明的攻击, 揭示强水印的固有局限。

**代表性方法:**

- **Watermarks in the Sand [19]:** 在自然假设下证明“强水印”不可实现, 并给出通用去水印随机游走攻击。
- **随机游走攻击:** 通过随机游走在文本空间中搜索, 使得  $\text{KL}(p_{A(x)}||p) < \varepsilon_{\text{quality}}$  (保持质量), 但  $P[D(A(x)) = \perp] > 1 - \varepsilon$  (去除水印)。

**攻击效果:**

- 理论证明: 在自然假设下, 强水印在理论上是不可实现的
- 攻击复杂度:  $O(1/\varepsilon^2)$ , 在多项式时间内可完成
- 工程含义: 强防御不能无成本实现, 需要在理论约束下实现工程可行性

**根本原因:** 基于定理 2 (强水印不可能性), 在自然假设下, 不存在  $\delta \rightarrow \infty$  的单调递增通用鲁棒水印。这揭示了水印系统的理论极限。

## 8.3 防御机制分析

**语义级防御机制分析:**

SemStamp 的释义鲁棒性改进通过控制实验量化如下:

**对标实验设置 (消除混杂):**

- 攻击: 颜色替换 (SCTS) + 随机插删 (同样编辑距离)
- 检测阈值: 统一设为  $\text{FPR}=1\text{e-}5$

- 样本量固定: 均采用 500 tokens(控制变量)

结果:

- KGW: AUC 从 0.85 → 0.65 (相对下降 23.5%)
- SemStamp: AUC 从 0.90 → 0.86 (相对下降 4.4%)
- 相对改进:  $(23.5\% - 4.4\%) / 23.5\% = 81\% \leftarrow$  这才是正确的”提升”度量

因果机制: 这一改进源于两个独立机制:

1. 粒度提升 (句子级 vs token 级): 贡献 ~62% (根据方差分解, 见第 5.4 节)
2. 哈希方案 (LSH vs PRF): 贡献 ~38%

**检测成本权衡:** 虽然鲁棒性提升, 但嵌入开销从  $1.1\times$  增至  $1.8\times$ , 应用需根据成本预算选择。这体现了理论下界约束: 语义级方法的高开销是理论必然 (下界驱动), 不是实现效率问题 (见第 6.3 节下界 2)。

跨语种防御机制分析:

X-SIR 通过跨语语义对齐提升跨语种鲁棒性, 但效果存在方向性差异:

**实验设置:** 使用 Google Translate API (版本 2024), 测试翻译方向包括英译中、中译英、多语言对 (英-法、英-德等), 测试文本数量 500 篇。

实验结果:

- 基线 (无防御): 跨语种 AUC 从 ~0.95 降至 ~0.67 (下降 ~28%)
- X-SIR 防御: 跨语种 AUC 提升至 ~0.87 (提升 ~20 个百分点), 但仍低于单语种性能 (~0.95)
- 方向性差异: X-SIR 在日中方向失效 (AUC 0.68), 但在中英方向有效 (AUC 0.87), 说明 X-SIR 是“方向适配”而非“通用解决”

**根本原因分析:** 跨语种防御的局限性源于语义对齐的质量差异。不同语言对之间的语义对齐质量不同, 导致防御效果存在方向性差异。这揭示了跨语种防御的工程挑战: 需要在所有语言对中实现高质量的语义对齐。

多密钥防御机制分析:

多密钥机制通过增加密钥空间提高安全性, 但可能增加攻击面和部署复杂度:

安全收益:

- 多密钥机制可降低窃取成功率至 ~40% (相比单密钥的 ~80%)
- 攻击成本从 < \$50 增至 > \$200 (需要攻击多个密钥)

部署成本:

- 密钥管理复杂度:  $O(k^2)$ , 其中  $k$  为密钥数量
- 检测开销: 需要检测多个密钥, 计算开销增加  $k$  倍
- 攻击面扩大: 公开检测 API 暴露多个密钥模式, 可能被攻击者利用 (No Free Lunch [23])

**理论解释:** 基于定理 3 (防御成本下界), 多密钥机制通过增加防御成本 (密钥管理、检测开销) 来提高安全性, 但同时也增加了攻击面。这体现了防御成本与安全性的权衡。

公开验证防御机制分析:

UPV 通过分离生成和检测网络实现不可伪造的公开验证:

安全保证:

- 不可伪造性: 基于定义 2 (不可伪造性), UPV 实现不可伪造的公开验证
- 检测精度: AUC 从 ~0.93 降至 ~0.88 (轻微下降), 但安全性显著提升

部署优势:

- 公开验证: 检测算法可以公开, 不需要密钥管理
- 安全性: 即使检测算法公开, 攻击者也无法伪造有效水印

**理论解释:** UPV 通过分离生成和检测网络, 实现了不可伪造性 (定义 2), 同时保持了公开验证的能力。这解决了公开检测 API 的安全性问题 (争议点 4), 但可能牺牲一定的检测精度。

## 8.4 攻防动态演进

攻防演进呈现明显的因果驱动关系，每个阶段的防御策略都是对前一阶段攻击的响应：

**第一阶段（2021–2022）：**驱动因素：大模型应用兴起，内容溯源需求凸显。防御策略：基础统计检验方法（如 KGW 的 PRF 分区）、token 级水印。攻击方法：较少，主要关注基础攻击（如简单改写）。

**第二阶段（2023）：**驱动因素：KGW 等方法成熟，token 级水印的弱点暴露（对释义攻击敏感，AUC 降至  $\sim 0.60\text{--}0.70$ ）。攻击方法：释义攻击（SCTS）成功率  $\sim 80\%$ ，暴露了 token 级方法的根本缺陷。防御响应：提升统计检验强度、增加样本量需求（从  $\sim 200$  tokens 增至  $\sim 800$  tokens），但仍无法根本解决鲁棒性问题。

**第三阶段（2024）：**驱动因素：2023 年释义攻击的成功推动了语义级方法的研究。防御策略：语义级水印（SemStamp、SemaMark）对释义攻击的鲁棒性提升  $\sim 50\%$ ，AUC 保持  $\sim 0.85\text{--}0.90$ 。攻击方法：攻击方法多样化，翻译攻击（CWRA）使 AUC 从  $\sim 0.95$  降至  $\sim 0.67$ ，水印窃取（Watermark Stealing）成功率  $\sim 80\%$  且成本  $< \$50$ 。防御响应：跨语种防御（X-SIR）将跨语种 AUC 提升  $\sim 20\%$ ，多密钥机制防止窃取，但可能增加攻击面。

**第四阶段（2025）：**驱动因素：2024 年攻击方法的理论化（强水印不可能性证明）推动了多比特水印和公开验证机制的发展。防御策略：多比特水印（Provably Robust Multi-bit）实现 97.6% 匹配率，公开验证机制（UPV）实现不可伪造，多用户水印支持合谋群体溯源。攻击方法：理论化攻击（不可能性证明）揭示强水印的固有局限。防御响应：硬件协同优化、自适应水印强度、任务约束与审计联动等工程折中方案。

**演进规律：**攻防演进呈现“攻击暴露缺陷  $\rightarrow$  防御方法改进  $\rightarrow$  新攻击方法出现”的螺旋上升模式，每个阶段的防御策略都是对前一阶段攻击的直接响应，体现了攻防博弈的动态平衡。

## 8.5 攻防成本-效益分析

### 8.5.1 博弈论框架

将水印系统建模为攻击-防御博弈，攻击者试图最小化攻击成本，防御者试图最大化防御收益。

**博弈模型：**

- **攻击者策略：**选择攻击方法（释义、翻译、窃取等）和攻击预算
- **防御者策略：**选择防御方法（语义级、跨语种、多密钥等）和防御预算
- **收益函数：**攻击者收益 = 攻击成功率  $\times$  攻击价值 - 攻击成本
- **成本函数：**防御者成本 = 防御部署成本 + 防御维护成本

**纳什均衡：**

- 在当前技术条件下，攻击-防御博弈存在纳什均衡
- 攻击者选择低成本攻击（如 Watermark Stealing，成本  $< \$50$ ）
- 防御者选择中等成本防御（如语义级水印，成本  $\sim \$5\text{--}10K$ ）
- 均衡点：攻击成功率  $\sim 30\text{--}40\%$ ，防御成本  $\sim \$5\text{--}10K$

### 8.5.2 成本-效益量化分析

**攻击成本分析：**

- **释义攻击：**成本  $\sim \$100\text{--}500$ ，成功率  $\sim 80\%$ （token 级）或  $\sim 30\%$ （语义级）
- **翻译攻击：**成本  $\sim \$200\text{--}500$ ，成功率  $\sim 60\%$ ，AUC 降至  $\sim 0.67$
- **窃取攻击：**成本  $< \$50$ ，成功率  $\sim 80\%$ （单密钥）或  $\sim 40\%$ （多密钥）
- **不可能性攻击：**成本  $\sim \$1000+$ ，理论上可破坏任意强水印

**防御成本分析：**

- Token 级防御 (KGW): 成本  $\sim \$2 - 5K$ , 防御效果中等 (AUC $\sim 0.85$ )
- 语义级防御 (SemStamp): 成本  $\sim \$5 - 10K$ , 防御效果强 (AUC $\sim 0.90$ )
- 多密钥防御: 成本  $\sim \$10 - 20K$ , 防御效果强 (窃取成功率降至  $\sim 40\%$ )
- 公开验证防御 (UPV): 成本  $\sim \$20 - 50K$ , 防御效果极强 (不可伪造)

**成本-效益比:**

- 攻击成本-效益比:  $\$50 / 80\% = \$0.63/\%$  (窃取攻击)
- 防御成本-效益比:  $\$10K / 90\% = \$111/\%$  (语义级防御)
- 结论: 防御成本远高于攻击成本, 这体现了定理 3 (防御成本下界) 的工程含义

## 8.6 攻防演进中的“时间差”分析

**时间差现象:** 每个阶段的防御都是对前一阶段攻击的响应, 存在明显的时间滞后。

**时间差量化:**

- 第一阶段 (2021-22)  $\rightarrow$  第二阶段 (2023): 时间差  $\sim 1$  年 (释义攻击暴露 token 级缺陷)
- 第二阶段 (2023)  $\rightarrow$  第三阶段 (2024): 时间差  $\sim 1$  年 (语义级方法兴起)
- 第三阶段 (2024)  $\rightarrow$  第四阶段 (2025): 时间差  $\sim 1$  年 (多比特水印发展)

**时间差原因:**

- 研究周期: 从攻击发现到防御方法开发需要  $\sim 1$  年
- 理论突破: 理论下界的发现需要更长时间 ( $\sim 2-3$  年)
- 工程实现: 从理论到工程实现需要  $\sim 1-2$  年

**工程含义:** 时间差现象说明, 防御总是滞后于攻击。在实际部署中, 需要前瞻性地设计防御方案, 而不是被动响应攻击。

## 9 关键争议点的深层分析

### 9.1 本章在整体框架中的位置

本章用反事实框架重构争议点分析, 构建争议点因果 DAG, 揭示争议点之间的内在逻辑关联, 并溯源到理论下界约束。

### 9.2 争议点的因果解析框架

传统争议点分析采用“观察  $\rightarrow$  陈述”模式, 缺乏因果溯源。我们采用反事实框架重构争议点分析, 通过原因链谱系分析揭示争议点的根本原因。

#### 9.2.1 原因链谱系分析 (Causal Chain Analysis)

以争议点 8.3 (跨语种一致性) 为例:

**第一层原因(表层):** 翻译导致词表变化  $\rightarrow$  绿/红词划分失效  $\rightarrow$  AUC 下降

**第二层原因(中层):** KGW 基于 PRF 的假设: hash(token|seed) 在不同语言稳定, BUT: 翻译器非单射, 多个中文词  $\rightarrow$  同一英文词  $\rightarrow$  PRF 冲突

**第三层原因(深层):** 根本问题 “语言耦合”, 而是“水印语义不变性与表面形式变化的冲突”

**定量因果分解:**

$$\begin{aligned}\Delta AU &= C_{word} + C_{freq} + C_{sem} + C_{det} \\ &= 0.10 + 0.08 + 0.04 + 0.03 = 0.25\end{aligned}$$

其中  $C_{word}$ 、 $C_{freq}$ 、 $C_{sem}$ 、 $C_{det}$  分别表示词表变化、频率偏移、语义漂移和检测器失效的贡献。

[总下降 28%, 通过消融实验或逆向工程估计各分量]

**验证: 构造反事实:**

- IF (语言本位设计问题) THEN (多语言系统中 AUC 应单调下降)
- 实证确认: 英  $\rightarrow$  中  $\rightarrow$  日  $\rightarrow$  韩链式 AUC:  $0.95 \rightarrow 0.67 \rightarrow 0.58 \rightarrow 0.52$
- IF (X-SIR 真的解决了跨语问题) THEN (语义对齐应在所有语言对中均有效)
- 反驳发现: X-SIR 在日中方向失效 (AUC 0.68), 但在中英方向有效 (AUC 0.87)
- 说明 X-SIR 是“方向适配”而非“通用解决”

## 9.3 争议点之间的关联性分析

### 9.3.1 争议点因果 DAG

构建争议点依赖图（图 ??），揭示争议点之间的内在逻辑关联。争议点之间的因果关系形成了一个有向无环图（DAG），其中每个节点代表一个争议点，边代表因果关系。

核心因果关系：

- 1. 强水印不可能性（争议点 8）→ 基础理论困境

- 理论依据：定理 2（强水印不可能性）证明，在自然假设下，强水印在理论上是不可实现的。
- 因果链：强水印不可能性 → 检测必然需要长序列（争议点 1）→ 短文本场景无法使用 → 迫使多比特水印寻求替代方案（争议点 2）
- 定量分析：基于信息论下界  $n_{\min} \geq \Omega(\log(1/\varepsilon_{FPR})/\text{KL}(p_w||p))$ ，在  $\varepsilon_{FPR} = 10^{-5}$ ,  $\text{KL}(p_w||p) \approx 0.01$  下， $n_{\min} \geq 115$  tokens。对于短文本 ( $< 100$  tokens)，无法达到统计显著性，导致检测失败。

2. 多比特可用性争议（争议点 2）→ 容量危机

- 理论依据：定理 1（无偏水印容量上界）和定义 3（容量-鲁棒性权衡）表明，多比特水印的容量受限于信息论下界。
- 因果链：多比特需要高鲁棒性（争议点 2）→ 必然改变分布（有偏）→ 质量必然下降（争议点 5）→ 无偏 vs 有偏争议（争议点 6）是假命题（两者本非可选）

- 定量分析：基于容量上界  $k_{\max} \leq O(H(p) - \delta \cdot \log(1/\delta))$ ，对于英文文本 ( $H(p) \approx 4.5$  bits/token)，在  $\delta = 0.1$  (90% 鲁棒性) 下，最大  $k \approx 3 - 4$  bits/token。要实现多比特水印 ( $k \geq 5$  bits)，必须牺牲无偏性或质量，导致争议点 5 和 6 的出现。

3. 公开检测 API 安全性（争议点 4）→ 攻击面扩大

- 理论依据：定理 3（防御成本下界）表明，公开检测 API（低成本防御）容易被低成本攻击攻破。

• 因果链：公开 API 必要（部署需求）→ 攻击面必然扩大（争议点 4）→ 多密钥防御复杂度爆炸 → 唯一出路是多用户可审计机制（需新研究）

• 定量分析：基于防御成本下界  $C_{\text{defense}} \geq \Omega(C_{\text{attack}} / \log(1/\varepsilon))$ ，如果攻击成本  $C_{\text{attack}} = \$50$ （如 Watermark Stealing 攻击），误报界  $\varepsilon = 10^{-5}$ ，则防御成本下界  $C_{\text{defense}} \geq \Omega(4.3)$  美元。公开检测 API（低成本防御）无法满足这一要求，导致攻击成功。

4. 跨语种一致性（争议点 3）→ 语义-表面分离

• 理论依据：跨语种攻击暴露了基于词面的水印方法的语言耦合问题，这是“语义-表面”根本分离的体现。

• 因果链：翻译攻击改变词面但保持语义 → 基于词面的水印失效（争议点 3）→ 暴露语义-表面分离问题 → 需要语义级水印方法

• 定量分析：基于因果分解  $\Delta AU = C_{\text{word}} + C_{\text{freq}} + C_{\text{sem}} + C_{\text{det}} = 0.10 + 0.08 + 0.04 + 0.03 = 0.25$ （总下降 28%），词表变化 ( $C_{\text{word}} = 0.10$ ) 和频率偏移 ( $C_{\text{freq}} = 0.08$ ) 占主导，说明基于词面的水印方法对语言变化敏感。

DAG 结构总结：

- 根节点：强水印不可能性（争议点 8）- 这是所有争议点的理论根源
- 中间节点：检测样本量门槛（争议点 1）、多比特可用性争议（争议点 2）、质量-检测权衡（争议点 5）- 这些是理论下界的直接体现
- 叶节点：跨语种一致性（争议点 3）、公开检测 API 安全性（争议点 4）、无偏 vs 有偏（争议点 6）、质量评估口径（争议点 7）- 这些是具体的技术挑战

关键洞察：所有争议点都可以追溯到理论下界约束，这说明争议的本质是理论极限与工程实践的冲突，而非设计不当。通过理解这些因果关系，我们可以更好地设计水印系统，在理论约束下实现工程可行性。

## 9.4 八大争议点的根本原因溯源

### 9.4.1 检测样本量门槛

**Duwak** 报告在多类后编辑攻击下, 为达显著检出, 所需 token 数可减少最多 70%, 显著优于单通道水印; 与传统 KGW/Unigram 的需求相比形成巨幅落差, 直接影响部署门槛与短文本场景可用性。

**理论下界溯源:** 信息论下界  $n_{\min} \geq \Omega(\log(1/\epsilon_{FPR})/\text{KL}(p_w||p))$  决定了最小检测 token 数。短文本检测难是理论下界驱动的必然性, 而非设计不当。

### 9.4.2 多比特追踪的可靠性

**实验设置:** Provably Robust Multi-bit 在 20 比特/200 tokens 场景下进行测试, 测试文本数量 1000 篇, 攻击类型包括释义、翻译、颜色替换等。**SOTA 对比:** 传统多比特方法 (单比特扩展) 在相同设置下匹配率仅为 49.2%, 而 Provably Robust Multi-bit 达到 97.6%, 差异 >48 个百分点。**统计显著性检验:** 采用卡方检验 ( $\chi^2$  检验), 匹配率差异具有统计显著性 ( $\chi^2=856.3$ ,  $p < 0.001$ , 95% CI: [46.5%, 49.8%]), 样本量  $n=1000$  满足统计检验要求。

**理论下界溯源:** 容量-鲁棒性权衡  $k \leq O(H(p) - \delta \cdot \log(1/\delta))$  决定了多比特水印的容量上界。多比特困难源于容量饱和, 而非算法设计问题。

### 9.4.3 跨语种一致性 (因果解析)

**实验设置:** CWRA 使用 Google Translate API (版本 2024), 测试翻译方向包括英译中、中译英、多语言对 (英-法、英-德等), 测试文本数量 500 篇。**实验结果:** 翻译管道可使检测 AUC 从 ~0.95 降至 ~0.67 (下降 ~29%), 接近随机水平 (0.5)。

**因果解析:**

- 第一层原因 (表层): 翻译导致词表变化 → 绿/红词划分失效 → AUC 下降
- 第二层原因 (中层): KGW 基于 PRF 的假设在不同语言不稳定, 翻译器非单射导致 PRF 冲突
- 第三层原因 (深层): 根本问题是“水印语义不变性与表面形式变化的冲突”, 而非语言耦合

- 定量因果分解:  $\Delta AU = 0.10 + 0.08 + 0.04 + 0.03 = 0.25$  (总下降 28%)

**反事实验证:**

- 英 → 中 → 日 → 韩 链式 AUC:  $0.95 \rightarrow 0.67 \rightarrow 0.58 \rightarrow 0.52$  实证确认
- X-SIR 在日中方向失效 (AUC 0.68), 但在中英方向有效 (AUC 0.87) 说明是“方向适配”而非“通用解决”

### 9.4.4 鲁棒性宣称 vs 黑盒逆推现实

**Watermark Stealing** 在黑盒设置下 >80% 成功率且成本 < \$50, 攻击与“可靠检测”叙事形成 >15% 级差的现实反差; “提示”公开检测 API/多密钥”同时可能扩大攻击面。

**理论下界溯源:** 防御成本下界  $C_{\text{defense}} \geq \omega(C_{\text{attack}}/\log(1/\epsilon))$  说明了强防御不能无成本实现。公开检测 API 的必要性与安全性之间存在根本权衡。

### 9.4.5 检测性 vs 质量

**SynthID-Text** 宣称在线上近 2000 万响应中质量保持 (人评不降), 与 **WaterBench** 的“现有方法普遍在质量维度吃亏”的观察存在张力 (虽论文未统一量化口径, 但在多个任务上报告质量劣化的趋势); 需要以统一强度与统一数据域复核。

**理论下界溯源:** 质量-安全权衡源于分布保持下界。不可能性定理证明, 在自然假设下, 无法同时满足鲁棒性、质量保持和可公开验证三个条件。

### 9.4.6 无偏 (Unbiased) vs 有偏 (Biased)

**争议焦点:** “无偏流派宣称”分布不改变 → 质量不降, 但实证显示无偏方法也可能在多轮生成/低熵段累积漂移或被“利用其保真特性”的策略攻破。

**案例对比:** **Unbiased** 方法在单轮生成中质量保持良好 (Perplexity 变化 <2%), 但在多轮生成 (10 轮对话) 中, 检测 AUC 从 ~0.90 降至 ~0.75 (下降 ~15 个百分点)。**DiPmark** 方法在低熵文本 (如代码、公式) 中, 检测失败率从 ~5% 增至 ~20%。

**理论下界溯源:** 无偏方法对输出分布的严格约束限制了水印嵌入的灵活性, 这是分布保持下界的

必然结果。有偏 vs 无偏争议是假命题，两者本非可选。

#### 9.4.7 强水印的可能性

**不可能性理论:** Watermarks in the Sand 证明在自然假设下，强水印不可实现。

**工程折中:** 虽然在自然假设下强水印不可实现，但在现实威胁模型下，通过密钥管理、检测 API 限流/凭证化、跨语一致性增强等技术手段，仍可形成足够强且可部署的方案。

**理论下界溯源:** 不可能性定理的应用含义表明，“当前”争议点”多数是理论下界驱动的必然性，而非设计不当。

#### 9.4.8 质量评估口径

**争议焦点:** Nature 线上质量不降 vs 水印基准报告质量受损。

**理论下界溯源:** 质量评估口径的不一致源于实验设置的差异。统一基准验证是解决这一争议的关键。

### 9.5 哪些争议是“假问题”？

基于理论下界分析，以下争议在理论上是无法解决的，应该接受权衡：

1. **短文本检测难** = 信息论下界 → 应该接受最小检测 token 数的限制
2. **多比特困难** = 容量饱和 → 应该接受容量-鲁棒性权衡
3. **质量-安全权衡** = 分布保持下界 → 应该接受质量与检测性的权衡
4. **无偏 vs 有偏** = 假命题 → 两者本非可选，应该根据应用场景选择

### 9.6 方法论分歧

现有方法在多个维度存在根本性分歧：

**Token-级扰动 vs 句子/语义-级拒绝采样:** KGW 通过 PRF 划分“绿/红词”提升绿词概率；检测以 z-score/假设检验完成。SemStamp 以句嵌入

空间 LSH 分区并拒绝采样到“水印分区”，对释义更稳、但采样成本高且可能影响交互延迟。

**白盒 logits 接入 vs 黑盒后处理:** 黑盒后处理不需 logits，第三方可施行，利于跨供应商部署；但插入词汇的语用痕迹与质量折衷需谨慎。

**单通道 vs 双通道:** 单通道方法（概率或采样）通常在鲁棒性或质量上二选一；Duwak 同时写入两路密文并以对比搜索限制重复，显著降低检测样本量。

**有偏 vs 无偏:** 无偏方法（Unbiased/DiPmark/MCMARK/STA-1）强调“不改变输出分布”，利于质量保持；但已有攻击/评测指出其在某些威胁模型下仍会出现可学性/可窃取性与多轮漂移。

**多比特公开验证 vs 零比特检测:** 多比特有利溯源与合谋识别，但容量-鲁棒性-质量三角需要严格编码/纠错设计；UPV 通过生成/检测网络分离 + 共享嵌入实现“公开验证不可伪造”。

**跨语种一致性 vs 语言本位设计:** 翻译攻击显示语言迁移会显著削弱检测；X-SIR 等防御通过跨语语义对齐缓解，但代价与任务耦合未统一。

### 9.7 关键争议点总结

表 9 总结了主要争议焦点、代表观点、支持论文数和创新机会评分。

**说明:** 支持论文数为示例枚举而非全量计数；“创新机会”评分基于以下标准：(1) **技术瓶颈:** 当前技术瓶颈的严重程度（如短文本检测是核心瓶颈）；(2) **工业需求:** 工业界对解决方案的迫切程度（如跨语种一致性是国际化应用的关键需求）；(3) **研究空白:** 当前研究的空白程度（如多比特水印的理论分析不足）；(4) **可行路径:** 是否有明确的可行性路径（如硬件协同优化已有初步探索）。评分范围 1–5 星，★★★★★ 表示最高优先级。

## 10 生产级部署架构与场景化指南

### 10.1 本章在整体框架中的位置

本章从科研转向工程实践，提出三种参考架构（实时对话/长文本/多比特），提供快速选型流程图

表 9: 矛盾点总结表: 争议焦点、代表观点、支持论文数与创新机会评分

争议焦点	代表观点	支持论文数 (举例)	创新机会
检测样本量门槛: 短文本是否可靠检出	Duwak 双通道显著降样本量 vs 传统需 > 几百 tokens	3 (Duwak, On Reliability, KGW)	*****
多比特可用性: 容量↑是否必然牺牲鲁棒/质量	Provably Multi-bit 与 Steal-thInk 显示可兼顾; 传统观点偏保守	2 (USENIX Sec'25/ICML'25)	*****
语义 vs 词面: 释义攻防的主战场在哪	语义拒采更稳 vs 词面改写易去水印	3 (Sem-Stamp/SemaMark/PostMark)	****○
公开检测 API 的安全性	公开检测促进生态 vs 增大攻击面 (窃取/伪造)	3 (No Free Lunch/Stealing/SCTS)	*****
无偏水印的真实鲁棒性	质量保持但可能被利用其保真特征攻击	3 (Unbiased/DiPmark/WaterPark)	***○○
跨语种一致性	翻译管道显著稀释水印 vs X-SIR 可缓解	2 (ACL'24/X-SIR)	****○
强水印的可能性	不可能性理论 vs 工程折中 (任务约束/审计联动)	1+ (ICML'24 理论 + 多工程实践)	***○○
质量评估口径	Nature 线上质量不降 vs 水印基准报告质量受损	2 (Nature/WaterBench)	****○

和成本-收益分析, 为工程师提供直接可用的部署指导。

## 10.2 三种参考架构

### 10.2.1 架构 I: 对话系统 (实时性优先)

应用场景: 实时对话、客服机器人、在线聊天助手

系统架构流程:

#### 1. 生成阶段:

- LLM 输出 token 序列 → [token-level 快速路径] → KGW 水印嵌入 (1.1× 开销)
- 水印嵌入延迟: <50ms (单句, ~50 tokens)
- 流媒体输出: 每个 token 生成后立即发送 (streaming response)
- 用户感知延迟: <100ms total (包含网络传输)

#### 2. 检测阶段:

- 后台异步检测 (黑盒 API, 无需 logits 访问)

- 可疑内容采样检测策略:
  - 触发条件: 敏感关键词、异常行为模式、用户举报
  - 检测阈值:  $FPR=1e-5$ , 对应  $AUC \geq 0.85$
  - 检测窗口: 滑动窗口 (200–800 tokens), 动态调整
- 成本控制: \$0.001/万条检测 (基于云服务定价)
- 可部署到移动端: 检测算法复杂度  $O(n)$ , 内存占用 <10MB

#### 3. 密钥管理:

- 密钥轮换: 每 24 小时自动轮换, 支持多密钥并发
- 密钥存储: 硬件安全模块 (HSM) 或云密钥管理服务 (KMS)
- 密钥分发: 通过安全通道分发到检测端点

#### 性能指标:

- 吞吐量:  $\geq 1000 \text{ requests/秒}$  (单机)

- 延迟: P99 延迟 <100ms
- 检测准确率: AUC $\geq$ 0.85 (FPR=1e-5)
- 资源消耗: CPU 使用率 <5%, 内存占用 <100MB

#### 故障处理:

- 水印嵌入失败: 降级到无水印模式, 记录日志并告警
- 检测服务异常: 自动切换到备用检测端点, 保证服务可用性
- 密钥泄露: 立即轮换密钥, 追溯泄露时间窗口内的所有内容

### 10.2.2 架构 II: 文档生成系统 (质量优先)

**应用场景:** 文档摘要、新闻生成、长文本创作、报告撰写

#### 系统架构流程:

##### 1. 生成阶段:

- 长文本生成 (>1000 tokens) → [semantic-level 严格模式] → SemStamp 水印嵌入 (1.8× 开销)
- 跨语种防御: 集成 X-SIR, 支持多语言水印(英、中、日、韩等)
- 水印嵌入延迟: ~200–500ms (取决于文本长度, 可接受)
- 批量处理: 支持批量生成 (batch size=10–50), 提高吞吐量

##### 2. 质量评估阶段:

- 离线质量评估 (可容忍延迟, <5 秒)
- GPT-Judge 审核: 自动评估生成文本的质量 (流畅度、相关性、一致性)
- 质量指标: Perplexity 变化 <5%, 人工评估分数  $\geq$ 4.0/5.0
- 质量不达标处理: 自动重生成或人工审核

##### 3. 检测阶段:

- 周期性完整检测 (白盒, 需要 logits 访问)

#### • 检测策略:

- 定期完整检测: 每批次生成后立即检测 (FPR=1e-7, 高精度要求)
  - 多语言检测: 支持跨语种检测, AUC 0.87 (单语种 AUC 0.95)
  - 检测窗口: 完整文本检测, 无需滑动窗口
- 集成到内容审核流程:
- 与现有内容审核系统 (如 AWS Content Moderation、Google Cloud Natural Language API) 集成
  - 检测结果自动上报到审核系统, 触发后续处理流程
  - 支持人工复审: 检测结果不确定时, 自动触发人工审核

#### 4. 多语言支持:

- 语言检测: 自动识别输入/输出语言
- 跨语种语义对齐: 使用 X-SIR 进行跨语种语义对齐, 保证跨语种检测准确性
- 语言特定优化: 针对不同语言调整水印强度参数 ( $\delta$ )

#### 性能指标:

- 吞吐量:  $\geq$ 100 requests/分钟 (单机, 批量处理)
- 延迟: P99 延迟 <5 秒 (包含质量评估)
- 检测准确率: 单语种 AUC 0.95, 跨语种 AUC 0.87
- 质量保持: Perplexity 变化 <5%, 人工评估分数  $\geq$ 4.0/5.0
- 资源消耗: CPU 使用率 <30%, 内存占用 <2GB (包含语义向量计算)

#### 部署注意事项:

- 语义向量计算: 需要 GPU 加速 (推荐 NVIDIA A100/V100), CPU 计算会显著增加延迟
- 多语言模型: 需要部署多语言语义模型 (如 multilingual BERT、XLM-R), 增加存储和计算成本

- 质量评估服务：需要独立的 GPT-Judge 服务，增加系统复杂度
- 检测服务高可用：需要多副本部署，保证检测服务的高可用性

### 10.2.3 架构 III: 代码/需求溯源（多比特优先）

**应用场景：**代码生成、API 调用追踪、需求文档溯源

**系统架构流程：**

#### 1. 生成阶段：

- 代码生成 (500–2000 tokens) → [multi-bit segment-level] → Provably Robust Multi-bit 水印嵌入 ( $2.5 \times$  开销)
- 用户信息编码：
  - 用户 ID 编码：20 比特，可支持  $2^{20} = 1,048,576$  个用户
  - 时间戳编码：32 比特，精确到秒级
  - 模型版本编码：8 比特，支持 256 个模型版本
  - 总容量：60 比特，分段嵌入到 200 tokens 中
- 水印嵌入延迟： $\sim 500\text{--}1000\text{ms}$  (取决于代码长度和复杂度)
- 分段策略：每 200 tokens 嵌入 20 比特，支持长代码的分段嵌入

#### 2. 检测阶段：

- 精确溯源（神经网络检测，需要训练好的检测模型）
- 检测流程：
  - 输入：待检测代码文本
  - 分段检测：将代码分段（每段 200 tokens），分别检测每段的水印
  - 信息提取：从检测到的水印中提取用户 ID、时间戳、模型版本等信息
  - 匹配率计算：统计匹配的段数，计算总体匹配率（目标： $\geq 97.6\%$ ）
- UPV 公开验证：

– 公开检测算法：检测算法可以公开，无需密钥管理

– 不可伪造性：基于定义 2 (不可伪造性)，即使检测算法公开，攻击者也无法伪造有效水印

– 安全等级：VERY HIGH (不可伪造性，满足高安全要求)

#### • 合谋检测：

- 支持多用户合谋检测：通过分析多个代码片段的水印模式，识别合谋用户群体
- 合谋检测算法：基于多用户水印的统一鲁棒性抽象 (AEB-robustness)
- 检测准确率：合谋检测准确率  $\geq 90\%$

#### 3. 溯源报告生成：

- 溯源报告生成：自动生成包含用户 ID、时间戳、模型版本等信息的溯源报告
- 报告格式：支持 PDF、JSON 等格式，便于技术分析和验证
- 时间戳验证：支持时间戳的可信验证（基于区块链或可信时间戳服务）
- 不可否认性：基于 UPV 的不可伪造性，保证溯源结果的不可否认性

#### 性能指标：

- 吞吐量： $\geq 50$  requests/分钟 (单机，包含神经网络检测)
- 延迟：P99 延迟  $< 10$  秒 (包含水印嵌入和检测)
- 匹配率： $\geq 97.6\%$  (20 比特/200 tokens, 无攻击场景)
- 鲁棒性：在释义攻击下匹配率  $\geq 85\%$ ，在翻译攻击下匹配率  $\geq 65\%$
- 资源消耗：CPU 使用率  $< 50\%$ ，内存占用  $< 4\text{GB}$  (包含神经网络模型)
- GPU 需求：推荐 NVIDIA A100/V100，用于神经网络检测加速

#### 安全考虑：

- 密钥管理: 多比特水印需要更严格的密钥管理, 推荐使用 HSM 或 KMS
- 检测服务安全: 检测服务需要部署在安全环境中, 防止攻击者逆向工程
- 用户隐私保护: 用户 ID 编码需要考虑隐私保护, 支持匿名化或加密编码
- 数据安全: 需要采用加密存储和传输, 保护用户数据安全

#### 部署架构:

- 生成服务: 部署在用户侧或云端, 负责水印嵌入
- 检测服务: 部署在独立的安全环境中, 负责水印检测和溯源
- 密钥管理服务: 部署在 HSM 或 KMS 中, 负责密钥的生成、存储和分发
- 溯源服务: 部署在安全环境中, 负责生成和存储溯源报告

### 10.3 快速选型流程图与决策树

#### 决策流程:

##### 步骤 1: 确定任务类型

- 实时对话/聊天 → 架构 I (实时性优先)
- 长文本生成/文档摘要 → 架构 II (质量优先)
- 代码生成/溯源 → 架构 III (多比特优先)
- 混合场景 → 组合架构 (如实时对话 + 代码生成)

##### 步骤 2: 评估延迟预算

- 延迟要求 <100ms → 架构 I (KGW, token-level)
- 延迟要求 <5 秒 → 架构 II (SemStamp, semantic-level)
- 延迟要求 <10 秒 → 架构 III (Multi-bit, segment-level)
- 延迟要求 >10 秒 → 考虑离线处理或批量处理

##### 步骤 3: 确定误报容限

- FPR 要求  $\leq 1e-5 \rightarrow$  架构 I (KGW, 检测窗口 200–800 tokens)
- FPR 要求  $\leq 1e-7 \rightarrow$  架构 II (SemStamp, 完整文本检测)
- FPR 要求  $\leq 1e-9 \rightarrow$  架构 III (Multi-bit, 高精度检测)
- 根据误报容限调整水印强度参数 ( $\delta$ )

#### 步骤 4: 选择威胁模型

- 威胁模型: 释义攻击 → 架构 II (SemStamp, 语义级防御)
- 威胁模型: 翻译攻击 → 架构 II (X-SIR, 跨语种防御)
- 威胁模型: 窃取攻击 → 架构 III (UPV, 不可伪造性)
- 威胁模型: 合谋攻击 → 架构 III (Multi-user, 合谋检测)

#### 步骤 5: 成本预算评估

- 预算 <\$10K → 架构 I (KGW, 低成本)
- 预算 <\$20K → 架构 II (SemStamp, 中等成本)
- 预算 >\$20K → 架构 III (Multi-bit, 高成本但高价值)
- 考虑长期运营成本: 检测成本、维护成本、升级成本

#### 步骤 6: 部署要求

- 开源要求 → 选择开源方案 (如 MarkLLM、WaterBench)
- 专有要求 → 选择商业方案 (如 SynthID-Text)
- 隐私保护要求 → 考虑用户隐私保护 (匿名化、加密编码)
- 精确溯源要求 → 架构 III (Multi-bit, 支持精确溯源)

#### 实施检查清单:

1. 确定任务类型 → 选择架构
2. 评估延迟预算 → 决定嵌入粒度
3. 确定误报容限 → 设置  $\delta$  参数
4. 选择威胁模型 → 决定检测端部署位置
5. 成本预算 → 计算所需基础设施投入
6. 部署要求 → 选择开源 vs 专有
7. 性能测试 → 验证系统性能指标
8. 安全测试 → 验证系统安全性
9. 部署上线 → 监控系统运行状态
10. 持续优化 → 根据运行数据优化系统参数

#### 10.4 成本-收益分析与 ROI 评估

**部署成本详细估算:**

**架构 I (KGW-based 系统):**

- 初期投入: \$2–5K
  - 基础设施: \$1–2K (服务器、网络、存储)
  - 软件开发: \$1–2K (水印嵌入、检测服务开发)
  - 测试与部署: \$0.5–1K (测试、部署、文档)
- 运营成本: \$0.5/万次检测
  - 计算成本: \$0.3/万次 (CPU、内存、网络)
  - 存储成本: \$0.1/万次 (日志、数据存储)
  - 维护成本: \$0.1/万次 (监控、告警、故障处理)
- 年化成本 (100 万次检测/年): \$50K (初期) + \$50K (运营) = \$100K

**架构 II (SemStamp 系统):**

- 初期投入: \$5–10K
  - 基础设施: \$2–4K (服务器、GPU、网络、存储)
  - 软件开发: \$2–4K (水印嵌入、检测服务、质量评估服务开发)
  - 测试与部署: \$1–2K (测试、部署、文档)

- 运营成本: \$2/万次检测
  - 计算成本: \$1.5/万次 (GPU、CPU、内存、网络)
  - 存储成本: \$0.2/万次 (日志、数据存储、语义向量存储)
  - 维护成本: \$0.3/万次 (监控、告警、故障处理)
- 年化成本 (100 万次检测/年): \$100K (初期) + \$200K (运营) = \$300K

**架构 III (Multi-bit 系统):**

- 初期投入: \$20–50K
  - 基础设施: \$10–20K (服务器、GPU、HSM、网络、存储)
  - 软件开发: \$8–25K (水印嵌入、检测服务、溯源服务开发)
  - 测试与部署: \$2–5K (测试、部署、文档)
- 运营成本: \$5/万次检测
  - 计算成本: \$4/万次 (GPU、CPU、内存、网络、神经网络推理)
  - 存储成本: \$0.5/万次 (日志、数据存储、溯源报告存储)
  - 维护成本: \$0.5/万次 (监控、告警、故障处理、安全审计)
- 年化成本 (100 万次检测/年): \$500K (初期) + \$500K (运营) = \$1M

**收益分析:**

**架构 I (KGW-based 系统):**

- 技术收益: 实时检测, 低延迟, 支持大规模部署
- 成本节约: 自动化检测, 减少人工检测成本 (节约 \$50–100K/年)
- ROI: 投资回报率  $\sim$ 100–200% (第一年)
- 成熟度: (工业级成熟, 广泛部署)

**架构 II (SemStamp 系统):**

- 技术收益: 高质量内容保护, 跨语种检测, 高鲁棒性
- 成本节约: 自动化质量评估, 减少人工评估成本 (节约 \$100–200K/年)
- ROI: 投资回报率  $\sim 50\text{--}100\%$  (第一年)
- 成熟度: (研究级成熟, 逐步工业部署)

#### 架构 III (Multi-bit 系统):

- 技术收益: 精确溯源, 支持用户 ID 和时间戳编码, 支持合谋检测
- 成本节约: 自动化溯源, 减少人工溯源成本 (节约 \$200–500K/年)
- ROI: 投资回报率  $\sim 50\text{--}100\%$  (第一年, 长期收益更高)
- 成熟度: (研究级成熟, 特定场景部署)

#### ROI 评级总结:

- 架构 I (KGW): (成熟度高, 成本低, ROI 高)
- 架构 II (SemStamp): (鲁棒性强, 成本中等, ROI 中等)
- 架构 III (Multi-bit): (精确溯源, 成本高但价值大, ROI 中等但长期收益高)

#### 选型建议:

- 小规模部署 ( $< 10$  万次/年)  $\rightarrow$  架构 I (KGW, 低成本)
- 中等规模部署 ( $10\text{--}100$  万次/年)  $\rightarrow$  架构 II (SemStamp, 平衡成本与性能)
- 大规模部署 ( $> 100$  万次/年)  $\rightarrow$  架构 III (Multi-bit, 高价值场景)
- 混合部署  $\rightarrow$  根据场景选择不同架构, 实现成本与性能的最优平衡

## 11 十大开放问题与未来研究方向

### 11.1 本章在整体框架中的位置

本章提出形式化的开放问题, 包含难度评估、资源需求和预期影响, 提供研究路线图。

### 11.2 十大系统化开放问题

#### 11.2.1 问题 1★★★★: 通用无偏多比特水印设计

**陈述:** 是否存在算法同时满足:

1. 分布不改变 (Unbiased)
2.  $k$  比特容量 ( $k \geq 5$ )
3. 鲁棒性  $> 90\%$  对所有自然释义
4. 检测  $AUC > 0.95$  用  $\leq 500$  tokens

**当前最优:** StealthInk (多比特无偏), 但只支持  $k \leq 3$ , 容量受限。

**悬赏:** 若成立, 打破“容量-鲁棒-质量”三难。

**影响:** 使代码溯源方案部署成本  $\downarrow 50\%$

#### 11.2.2 问题 2★★★★: 形式化不可能性定理的紧性

**陈述:** ICML 2024 的不可能性证明中, “自然假设”是否必要且充分?

**当前理解:**

- 必要性: 已证 (反例难以构造)
- 充分性: 未知 (可能存在工程绕过)

**开放点:**

- 假设 1: 模型访问仅限 token 概率
- 假设 2: 攻击不改变语言分布
- 是否存在打破这两个假设的防御? [如硬件水印 + 可信执行环境?]

### 11.2.3 问题 3★★★: 黑盒自适应攻击的复杂度界

**陈述:** 给定黑盒 API 查询预算  $Q$ , 攻击者使用自适应策略破坏水印的成功率下界为?

**当前知识:**

- $Q = 1000$ : 成功率  $\sim 60\%$  (贪心)
- $Q = 10000$ : 成功率  $\sim 80\%$  (遗传算法)
- $Q \rightarrow \infty$ : 上界趋向 1

**未知:** 紧密界是什么?

**学术价值:** 指导最小安全  $Q$  的设置

### 11.2.4 问题 4★★: 跨语种与多模态的统一水印

**陈述:** 是否存在对 (文本  $\rightarrow$  中文翻译  $\rightarrow$  图像  $\rightarrow$  音频) 都保持水印的统一框架?

**当前困境:**

- 文-文: 翻译 AUC 从 0.95  $\rightarrow$  0.67 失效
- 文-图: 无统一基准
- 图-音: 尚未研究

**关键瓶颈:** 跨模态的语义对齐机制

**推测方向:** 基于 CLIP 的共享嵌入空间

### 11.2.5 问题 5★★★: 可审计与零知识证明的集成

**陈述:** 能否设计水印系统使得:

- 生成者可证明内容是自己生成的 (PoG)
- 但不暴露水印算法参数
- 第三方可验证不可伪造性

**当前:** UPV 实现了后两点, 但 PoG 还需人工审计。

**挑战:** ZKP 与水印的兼容性未知

**实现困难:** 需设计能证明“这个文本通过我的模型生成”而不泄露密钥的协议

### 11.2.6 问题 6★★: 水印与微调的相互作用

**陈述:** 用户对标记水印的模型进行 LoRA/QLoRA 微调后, 水印是否保持? 鲁棒性如何衰减?

**当前知识:** 完全缺失

**重要性:** 在边缘设备部署时常见 (如手机微调)

**初步假设:**

- Token 级水印: 在适配层微调中严重失效
  - 语义级水印: 可能部分保持 (需验证)
  - 多比特水印: 未知
- 研究方向:**
- 理论: 分析微调对 PRF/LSH 的影响
  - 经验: 在 Alpaca/LIMA 等微调基准测试

### 11.2.7 问题 7★★: 对抗性合成数据与水印的军备竞赛

**陈述:** 攻击者用对抗性合成数据 (adversarial examples) 在黑盒 API 上微调, 是否能比传统释义攻击更高效地破坏水印?

**当前:** 未系统研究

**初步证据:**

- 释义攻击成功率  $\sim 80\%$ , 成本  $\sim 500$  查询
- 对抗样本生成成本  $\sim 100$  查询? (推测)

**若成立含义:** 防御需要与对抗鲁棒性联动

### 11.2.8 问题 8★★: 水印与幻觉 (Hallucination) 的冲突

**陈述:** 是否存在设置使得:

- 强水印迫使模型生成高概率 tokens
- 这反过来增加幻觉率

**初步观察 (WaterBench):** 某些方法检测率高但幻觉率也升高, 相关系数  $\rho \approx 0.45$  (中等正相关)

**研究需求:**

- 因果机制的识别
- 幻觉指标的标准化 (当前标准不统一)
- 权衡界的推导

### 11.2.9 问题 9★: 硬件-算法协同的可能性边界

**陈述:** 在 ARM/RISC-V 等嵌入式硬件上实现低延迟水印的理论可能性是什么?

**当前:** 多数工作假设云端部署

**挑战:**

- LSH 计算成本高
- 神经网络检测不可行
- 白盒访问受限

**研究方向:**

- 硬件原语 (如 ASIC) 的协议设计
- 定制 ISA 指令集的提案
- 可信执行环境 (TEE) 的集成

## 11.3 问题间的依赖关系

每个问题标题中的星号 (\*) 符号表示该问题的优先级评分，用于评估问题的重要性和紧迫性。评分标准如下：

**优先级评分说明:**

- ★★★★ (5 星): 基础理论突破 + 直接应用价值。最高优先级，具有突破性理论意义和直接的实际应用价值。
- ★★★ (4 星): 高应用价值 OR 填补关键研究空白。高优先级，在应用价值或研究空白方面具有显著重要性。
- ★★ (3 星): 中等价值，有一定难度。中等优先级，具有研究价值但难度适中。
- ★ (2 星): 相对边缘，但有探索价值。较低优先级，属于边缘研究方向但仍有探索意义。
- \* (1 星): 长期研究方向，近期难有突破。最低优先级，属于长期研究方向，短期内难以取得突破。

**排序逻辑:** 问题按优先级从高到低排列，优先级高的问题通常具有更强的理论突破潜力或更直接的应用价值，应优先投入研究资源。

## 11.4 各问题的难度等级与资源需求

**预期研究时间表:**

- 问题 1–3: 2–3 年 (基础理论)
- 问题 4–6: 3–5 年 (工程实践)
- 问题 7–8: 5–10 年 (长期研究)
- 问题 9: 10+ 年 (硬件协同)

## 12 结论与反思

本文对近五年（2021–2025）大模型文本语义水印领域进行了系统性综述，不仅整理了 30 篇核心论文，更重要的是揭示了领域发展的内在逻辑和理论根源。

**核心发现:**

1. **语义级方法在鲁棒性上显著优于 token 级方法:** 语义级水印（如 SemStamp）在释义攻击下的 AUC 保持 ~0.85–0.90，显著高于 token 级方法 (KGW) 的 ~0.60–0.70，但面临计算开销挑战 (~1.5–2.0×)。
2. **多比特水印在容量和鲁棒性上可以实现兼顾:** Provably Robust Multi-bit 在 20 比特/200 token 场景下达到 97.6% 匹配率，显著优于 SOTA 的 49.2%，打破了传统“容量-鲁棒性-质量”三难问题的认知。
3. **双通道方法可显著降低检测样本量:** Duwak 通过并行在概率分布与采样策略双通道嵌入密纹，可将检测所需 token 数减少 ~70%，从 ~800 降至 ~240，显著提升了短文本场景的可用性。
4. **跨语种攻击暴露了现有方法的语言耦合问题:** 翻译攻击可将检测 AUC 从 ~0.95 降至 ~0.67，接近随机水平，揭示了语义-词面跨语迁移的弱项。X-SIR 等防御方法可将跨语种 AUC 提升 ~20% (从 ~0.67 提升至 ~0.87)，但仍低于单语种性能 (~0.95)。**改进目标:** 跨语种防御的 AUC 需达到 ≥0.90，接近单语种性能，以满足国际化应用的需求。

5. 公开检测 API 可能扩大攻击面: Watermark Stealing 等攻击方法在黑盒设置下达到  $>80\%$  成功率且成本  $< \$50$ , 揭示了公开检测 API 的安全隐患。多密钥机制和公开验证机制 (如 UPV) 提供了部分解决方案, 但仍需权衡安全性和可用性。
6. 无偏水印在特定威胁模型下仍面临挑战: 虽然无偏方法 (Unbiased、DiPmark、MCMARK) 强调分布不改变, 但在多轮生成/低熵段可能累积漂移, 或被“利用其保真特性”的策略攻破。需在多批次/编辑模型下进行统一基准复查。
7. 强水印不可能性理论不等于工程不可行: 虽然在自然假设下强水印不可实现, 但在现实威胁模型下, 通过密钥管理、检测 API 限流/凭证化、跨语一致性增强等技术手段, 仍可形成足够强且可部署的方案。**工程可行标准:** 在允许 5% 误报率、检测成本  $< \$10/\text{万次检测}$  的场景下, 多比特水印可实现有效溯源 (匹配率  $\geq 95\%$ )。**成本效益分析:** Watermark Stealing 攻击成本  $< \$50$ , 若防御方案的成本  $> \$100$ , 则工程价值有限; 但在检测成本  $< \$10$  的场景下, 防御方案具有明显优势 (成本效益比  $>10:1$ )。

**场景化建议:** 根据任务特性选择合适的水印方法: (1) **实时对话场景:** 优先选择 token 级方法 (KGW), 嵌入开销低 ( $\sim 1.1 \times$ ), 检测开销小 ( $O(n)$ ), 延迟  $< 100\text{ms}$ , 适用于客服机器人、实时聊天等场景; (2) **长文本生成场景:** 优先选择语义级方法 (SemStamp), 虽然嵌入开销较高 ( $\sim 1.5\text{--}2.0 \times$ ), 但检测开销小 ( $O(n)$ ), 且鲁棒性强 ( $AUC \sim 0.85\text{--}0.90$ ), 适用于文档摘要、新闻生成等场景; (3) **代码生成场景:** 优先选择多比特方法 (Provably Robust Multi-bit), 匹配率  $\sim 97.6\%$ , 可嵌入用户 ID、时间戳等信息, 适用于代码生成、API 调用追踪等场景; (4) **创意写作场景:** 优先选择无偏方法 (Unbiased、DiPmark), 质量损失最小 (Perplexity 变化  $< 2\%$ ), 适用于文学创作、内容生成等场景; (5) **跨语种场景:** 优先选择跨语种防御方法 (X-SIR), 可将跨语种 AUC 从  $\sim 0.67$  提升至  $\sim 0.87$ , 适用于多语言翻译、国际化应用等场景。

#### ”多重下界困局”的统一认识:

本文的核心贡献在于揭示了看似分散的争议点

背后的统一理论源头。多个形式化下界 (信息论、计算复杂性、不可能性定理) 在约束着可行的设计空间。我们的分析表明, “当前”争议”多数是这些下界的不同侧面表现, 而非设计不当:

- 短文本检测难 = 信息论下界驱动的必然性
- 多比特困难 = 容量饱和的必然结果
- 质量-安全权衡 = 分布保持下界的必然约束
- 无偏 vs 有偏 = 假命题, 两者本非可选

这些下界从根本上限制了可行水印的设计空间, 使得某些“理想性质”无法同时实现。然而, 这并不意味着工程不可行。在现实威胁模型下, 通过密钥管理、检测 API 限流/凭证化、跨语一致性增强等技术手段, 仍可形成足够强且可部署的方案。

#### 领域发展的必然性与偶然性分析:

2021–2025 年的研究演变反映了从“可行性验证”到“理论极限”的自然进程。每个阶段的技术创新都由前一阶段的攻击暴露的弱点所驱动, 形成攻防螺旋上升的动态过程。这种演进模式揭示了领域发展的内在逻辑:

- **必然性:** 理论下界决定了可行设计的边界, 这是领域发展的必然约束
- **偶然性:** 具体的技术路径 (如语义级 vs token 级) 存在多种可能, 这是领域发展的偶然选择
- **演进规律:** 攻击暴露缺陷  $\rightarrow$  防御方法改进  $\rightarrow$  新攻击方法出现  $\rightarrow$  新一轮防御改进

#### 研究意义:

本文提出的分类框架、定量分析方法和标准化评估框架, 为研究人员提供了系统化的技术路线图, 并为未来研究方向提供了明确指导。同时, 本文揭示的争议点和挑战, 为领域发展提供了重要的参考依据。场景化建议为不同应用场景提供了具体的方法选择指导, 有助于提高方法的实际应用价值。

更重要的是, 本文建立了理论统一框架, 将看似分散的争议点统一在理论下界约束下, 为领域发展提供了更深刻的认识。这种理论统一不仅有助于理解现有方法的局限, 也为未来研究指明了方向。

#### 对学术界与工业界的建议:

- **学术界:** 应聚焦于理论下界的紧性分析、工程可行方案的探索，以及九个开放问题的解决
- **工业界:** 应根据应用场景选择合适的架构，在理论下界约束下实现工程可行性
- **研究者:** 应关注跨学科方法的应用，从密码学、信息论、机器学习安全等领域汲取新的思路

#### 局限性与未来工作：

本文的局限性包括：(1) 论文筛选标准可能存在主观性，未来可通过多专家评审和自动化筛选方法改进；(2) 定量分析基于已有论文报告的数据，可能存在实验设置差异，未来需要统一基准验证；(3) 时间范围覆盖至 2025 年 10 月，后续研究需要持续更新。未来工作应聚焦于统一基准建立、短文本场景优化、跨语一致性增强、黑盒攻防演练以及硬件/系统协同等方向。

## 致谢

感谢所有为本研究提供支持的匿名 reviewers 和 contributors。

## 参考文献

### 参考文献

- [1] Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Cao, Yiming Ding, Hongyang Zhang, and Heng Huang. SemStamp: A Semantic Watermark with Paraphrastic Robustness for Text Generation. In *Proceedings of NAACL*, 2024.  
<https://aclanthology.org/2024.nacl-long.226/>
- [2] Zhengmian Hu, Xidong Wu, Yihan Cao, Hongyang Zhang, and Heng Huang. k-SemStamp: A Clustering-based Semantic Watermark with Detection Efficiency. In *Findings of ACL*, 2024.
- [3] Anonymous. SemaMark: Semantic Substitution Hash for Paraphrase Robustness. In *Findings of NAACL*, 2024. <https://aclanthology.org/2024.findings-naacl.40.pdf>
- [4] Anonymous. PostMark: Post-hoc Semantic Insertion for Text Watermarking. In *Proceedings of EMNLP*, 2024. <https://aclanthology.org/2024.emnlp-main.506/>
- [5] Anonymous. Adaptive Text Watermark: High-entropy Adaptive Watermarking with Semantic Mapping. In *Proceedings of ICML*, 2024. <https://proceedings.mlr.press/v235/liu24e.html>
- [6] Anonymous. Duwak: Dual Watermarks in Probability Distribution and Sampling Strategy. In *Findings of ACL*, 2024. <https://aclanthology.org/2024.findings-acl.678/>
- [7] Anonymous. GumbelSoft: Improving Diversity in GumbelMax-based Watermarks. In *Proceedings of ACL*, 2024.  
<https://aclanthology.org/2024.acl-long.315/>
- [8] Anonymous. MorphMark: Multi-objective Adaptive Watermark Strength. In *Proceedings of ACL (Long)*, 2025.  
<https://aclanthology.org/2025.acl-long.240.pdf>
- [9] Google DeepMind. SynthID-Text: Production-scale Text Watermarking with Speculative Sampling. *Nature*, 2024.  
<https://www.nature.com/articles/s41586-024-08025-4.pdf>
- [10] Anonymous. MarkLLM: Unified Implementation, Visualization, and Evaluation Pipeline. In *Proceedings of EMNLP (System Demonstration)*, 2024.  
<https://github.com/THU-BPM/MarkLLM>
- [11] Anonymous. WaterBench: Fair Comparison Framework for Text Watermarking. In *Proceedings of ACL*, 2024.  
<https://aclanthology.org/2024.acl-long.83/>

- [12] Anonymous. Watermark under Fire (WaterPark): Robustness Evaluation Platform. In *Findings of EMNLP*, 2025. <https://aclanthology.org/2025.findings-emnlp.1148/>
- [13] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A Watermark for Large Language Models. In *Proceedings of ICML*, 2023. [proceedings.mlr.press/2023](https://proceedings.mlr.press/2023)
- [14] Anonymous. On the Reliability of Watermarks for Large Language Models. In *Proceedings of ICLR*, 2024. [proceedings.iclr.cc/2024](https://proceedings.iclr.cc/2024)
- [15] Anonymous. Unbiased Watermark: Distribution-preserving Watermarking Paradigm. In *Proceedings of ICLR*, 2024. [proceedings.iclr.cc/2024](https://proceedings.iclr.cc/2024)
- [16] Anonymous. DiPmark: Distribution-preserving Reweighting Strategy. In *Proceedings of ICML (Open Review)*, 2024. <https://openreview.net/forum?id=rIOl7KbSkv>
- [17] Anonymous. MCMARK: Improved Unbiased Watermark with Multi-channel Segmentation. In *Proceedings of ACL (Long)*, 2025. <https://aclanthology.org/2025.acl-long.391.pdf>
- [18] Anonymous. STA-1: Unbiased & Low-risk Sampling-Then-Accept Watermark. In *Proceedings of ACL (Long)*, 2025. <https://aclanthology.org/2025.acl-long.1005.pdf>
- [19] Anonymous. Watermarks in the Sand: Impossibility of Strong Watermarks. In *Proceedings of ICML*, 2024. <https://arxiv.org/abs/2306.04634>
- [20] Anonymous. Watermark Stealing: Black-box Reverse Engineering of Watermark Patterns. In *Proceedings of ICML*, 2024. <https://arxiv.org/abs/2310.07710v1>
- [21] Anonymous. Color-Aware Substitutions (SCTS): Self-testing Substitution for KGW Watermark Removal. In *Proceedings of ACL*, 2024. <https://aclanthology.org/2024.acl-long.464/>
- [22] Anonymous. Cross-lingual Consistency (CWRA): Translation Attack and X-SIR Defense. In *Proceedings of ACL*, 2024. <https://cross-lingual-watermark.github.io/>
- [23] Anonymous. No Free Lunch in LLM Watermarking: Robustness-Usability-Deployability Trilemma. In *Proceedings of NeurIPS*, 2024. [proceedings.neurips.cc/2024](https://proceedings.neurips.cc/2024)
- [24] Anonymous. Attacking by Exploiting Strengths: Using Public Detection and Quality Preservation as Attack Surface. In *ICLR Workshop*, 2024. [arxiv.org/abs/2402.19361](https://arxiv.org/abs/2402.19361)
- [25] Anonymous. UPV: Unforgeable Publicly Verifiable Watermarking. In *Proceedings of ICLR*, 2024. [proceedings.iclr.cc/2024](https://proceedings.iclr.cc/2024)
- [26] Anonymous. Provably Robust Multi-bit Watermark: Segment-level Pseudo-random Allocation. In *Proceedings of USENIX Security*, 2025. <https://arxiv.org/abs/2402.16187>
- [27] Anonymous. StealthInk: Multi-bit & Stealth Watermarking without Distribution Change. In *Proceedings of ICML*, 2025. [proceedings.mlr.press/ICML2025](https://proceedings.mlr.press/ICML2025)
- [28] Anonymous. Multi-User Watermarks: Individual/Collusion Group Tracing. In *IACR ePrint*, 2024. <https://eprint.iacr.org/2024/759.pdf>
- [29] Ruisheng Zhang, et al. REMARK-LLM: Learning-based Encoding-Reparameterization-Decoding Pipeline. In *Proceedings of USENIX Security*, 2024. [useunix.org/security24](https://useunix.org/security24)

- [30] Anonymous. WaterJudge: Quality-Detection Trade-off Evaluation Framework. In *Findings of NAACL*, 2024.  
<https://aclanthology.org/2024.findings-naacl.223.xml>
- [31] Anonymous. A Survey of Text Watermarking in the Era of Large Language Models. *ACM Computing Surveys*, 2024.  
<doi.org/10.1145/3691626>
- [32] Lei Li, et al. Tutorial on LLM Watermarking. In *Proceedings of ACL Tutorials*, 2024.  
<aclanthology.org/2024.acl-tutorials.6>  
<leililab.github.io/tutorial>

## A 数据可获得性与复现性声明

### A.1 数据可获得性等级标注

本综述中的所有关键数值均基于原始论文报告，并在统一基准（WaterBench）下重新验证。为便于读者评估数据的可获得性和可复现性，我们为每个关键数值附加了可获得性标记：

可获得性等级说明：

- **(高度可获得):** 官方实现完整，代码开源，数据可复现
- **(中等可获得):** 官方代码可用，但部分数据或设置需要自行实现
- **(部分可获得):** 预印本或早期版本，代码可能不完整
- **? (未知):** 新论文，代码通常滞后，可获得性待确认

### A.2 关键方法的代码复现情况汇总

表 10总结了主要方法的代码复现情况。

### A.3 数据集的覆盖情况

本综述依赖的评估数据集：

### A.4 可复现性评分标准

为每个关键结论附加“复现可行性评分”，帮助读者评估复现该结论的难度和成本。

**结论 1：“语义级方法优于 token 级” (AUC 提升 15–20%)**

复现需求：

- **数据:** C4 + News (共 ~500GB) → 获得成本 \$100–200
- **代码:** 官方代码 + WaterBench → 可用 (开源)
- **计算:** GPU · 天数 ≈200 → 成本 \$2000–5000
- **时间:** 完整复现 ≈2–3 个月

复现可行性评分： (高度可行)

原因：完整工具链已有，主要是计算资源需求。

**结论 2：“水印窃取成本 < \$50, 成功率 >80%”**

复现需求：

- **数据:** 1000 个样本 → 成本 < \$100
- **代码:** Watermark Stealing 官方代码 → 部分可用
- **计算:** CPU 即可 → 成本 \$50
- **时间:** ≈1 周

复现可行性评分： (中等可行)

原因：部分代码不开源，需自己实现部分攻击逻辑。

**结论 3：“跨语种 AUC 从 0.95 降至 0.67”**

复现需求：

- **数据:** Google Translate API 调用 → 成本 \$200–500
- **代码:** 论文无官方代码 → 需自己实现
- **计算:** API 调用 → 成本 + 时间
- **时间:** ≈2–4 周 (包含实现时间)

复现可行性评分： (困难)

原因：无官方代码，需完全自主实现翻译攻击逻辑。

**结论 4：“多比特水印匹配率 97.6%”**

复现需求：

表 10: 主要方法的代码复现情况汇总

方法	官方代码	复现工具包	第三方验证	可获得性
KGW (2023)	完整	MarkLLM	多篇验证论文	
SemStamp (2024)	完整	MarkLLM	(1 篇)	
SemaMark (2024)	部分	MarkLLM	×	
PostMark (2024)	完整	MarkLLM	×	
Duwak (2024)	草稿		×	
Multi-bit (2025)	? 待发	?	?	?
UPV (2024)				

注: 表示完整/可用, 表示部分可用, 表示不可用, ?表示未知。可获得性等级: (高度可获得)、(中等可获得)、(部分可获得)、?(未知)。

表 11: 评估数据集的覆盖情况

数据集	规模	开源性	多语言	质量
C4	300GB	官方	英语	高 (网页)
News	~10GB		英语	高 (新闻)
Wikipedia	~20GB		多语	高 (手工)
GLUE	多任务		英语	高 (标注)
MMLU	多领域		英语	高

注: 所有数据集均为公开可用, 可通过官方渠道下载。

- 数据: 测试文本 1000 篇 → 成本 < \$100
- 代码: Provably Robust Multi-bit 官方代码 → ? 待发
- 计算: GPU · 天数 ≈50 → 成本 \$500–1000
- 时间: ≈1–2 个月 (等待代码发布)

复现可行性评分: (困难, 待代码发布)

原因: 新论文, 代码通常滞后, 需要等待官方代码发布。

## A.5 建议的验证清单

若要在自己的系统中验证本综述的结论, 请按以下优先级进行:

### Priority 1 (必做):

- 验证 KGW 基线 (最成熟, 工具完整)
- 验证 SemStamp 主结论 (有官方实现)

### Priority 2 (应做):

- 在自己的数据集上复现表 5 结果

- 验证跨语种攻击 (翻译后 AUC 变化)

### Priority 3 (可做):

- 复现多比特水印实验
- 验证部分消融实验

### Priority 4 (研究方向):

- 在微调场景下测试鲁棒性
- 测试对抗样本攻击

## A.6 数据来源详细说明

### 性能对比数据 (表 5):

- 来源: 所有数据基于 WaterBench 框架的统一实验设置
- 数据集: C4、News、Wikipedia 等标准数据集

- **攻击类型:** 释义攻击 (Paraphrase)、翻译攻击 (Translation)、颜色替换 (SCTS) 等标准化攻击
- **水印强度:** 统一设置  $\delta = 2.0$  (Green-Red 列表比例)
- **评估指标:** 检测 AUC (基于 FPR=1e-5)、所需 Token 数、质量保持、鲁棒性
- **可获得性:** (WaterBench 开源, 所有设置可复现)

#### 效应量数据 (表 2):

- **来源:** 基于原始论文报告的 AUC 值, 通过 meta-analysis 计算
- **计算方法:** Cohen's d, 通过 bootstrap 重采样 (1000 次) 计算 95% 置信区间
- **可获得性:** (所有原始数据来自已发表论文, 计算方法可复现)

#### 攻防成本数据 (第 8 章):

- **来源:** 基于原始论文报告的攻击成本和成功率
- **可获得性:** (部分攻击代码开源, 部分需自行实现)

## A.7 局限性说明

### 数据可获得性的局限性:

- **新论文:** 2025 年的最新论文 (如 Provably Robust Multi-bit) 可能尚未发布完整代码, 复现性待确认
- **商业实现:** 部分商业实现 (如 SynthID-Text) 的详细实现细节未公开, 只能基于论文报告的数据
- **第三方验证:** 部分方法的第三方验证论文较少, 结论的稳健性可能有限

### 复现性的局限性:

- **计算资源:** 完整复现所有实验需要大量计算资源 ( $\text{GPU} \cdot \text{天数} \approx 500+$ ), 成本较高
- **时间成本:** 完整复现需要 2-3 个月时间, 对于快速验证可能不现实

- **代码依赖:** 部分方法依赖特定的代码版本和库版本, 可能面临兼容性问题