

Conference Paper Title*

*Note: Sub-titles are not captured for <https://ieeexplore.ieee.org> and should not be used

1st Given Name Surname

dept. name of organization (of Aff.)

name of organization (of Aff.)

City, Country

email address or ORCID

2nd Given Name Surname

dept. name of organization (of Aff.)

name of organization (of Aff.)

City, Country

email address or ORCID

3rd Given Name Surname

dept. name of organization (of Aff.)

name of organization (of Aff.)

City, Country

email address or ORCID

4th Given Name Surname

dept. name of organization (of Aff.)

name of organization (of Aff.)

City, Country

email address or ORCID

5th Given Name Surname

dept. name of organization (of Aff.)

name of organization (of Aff.)

City, Country

email address or ORCID

6th Given Name Surname

dept. name of organization (of Aff.)

name of organization (of Aff.)

City, Country

email address or ORCID

摘要—随着生成式扩散模型在图像内容生成领域的广泛应用，如何对生成内容进行有效的溯源和版权保护成为亟需解决的问题。本文提出了一种结合 EDICT 与 BDIA 的新型扩散模型水印方法：在前向生成阶段采用 EDICT 的耦合潜变量机制提升水印嵌入的精确性与鲁棒性，在逆向提取阶段引入 BDIA 高效积分近似方法，实现高效且精确的水印恢复。理论分析与实验结果均表明，该方法在保证生成性能无损的前提下，显著提升了水印的检测率、可追溯性和鲁棒性，并大幅降低了计算开销。本文为 AI 生成内容的可信标记与溯源提供了更高效、实用的新思路。

Index Terms—扩散模型，图像水印，性能无损，鲁棒性，内容溯源

I. 引言

近年来，扩散模型（Diffusion Models）在图像生成、内容创作等领域取得了突破性进展，成为生成式人工智能的重要基础。然而，随着 AI 生成内容（AIGC）的广泛传播，如何对生成图像进行有效的溯源和版权保护，防止内容滥用与伪造，成为学术界和产业界关注的焦点。图像水印作为数字内容溯源与版权保护的核心技术，面临着鲁棒性、容量、无损性等多重挑战。

Identify applicable funding agency here. If none, delete this.

现有扩散模型水印方法如 Gaussian Shading [1] 等，虽然能够在噪声潜空间嵌入水印，但在逆向提取时往往存在精度损失或计算开销过大的问题。近期，EDICT [2] 方法通过耦合潜变量实现了扩散过程的精确反演，但其逆向过程计算量大，效率有限。BDIA [3] 方法则提出了一种高效的双向积分近似机制，可在保证精度的同时大幅提升逆向推理效率。

为此，本文提出了一种结合 EDICT 与 BDIA 的新型扩散模型水印方法：在前向生成阶段采用 EDICT 的耦合潜变量机制提升水印嵌入的精确性与鲁棒性，在逆向提取阶段引入 BDIA 高效积分近似方法，实现高效且精确的水印恢复。该方法兼具精确性与高效性，理论与实验均证明其在水印检测率、可追溯性、鲁棒性和推理效率等方面优于现有扩散水印方法。

II. 相关工作

近年来，针对 AIGC 内容的溯源与版权保护，学术界提出了多种图像水印方法。传统水印方法主要包括基于空域和频域的嵌入方式，如在像素空间或 DCT/DFT 等频域对图像进行微扰。这些方法在静态图像中具有一定的鲁棒性，但在扩散模型生成流程中，嵌入的水印极

易被扩散过程中的噪声破坏，导致检测率和可追溯性大幅下降。

为解决上述问题，部分研究尝试将水印嵌入扩散模型的生成流程中。例如，Gaussian Shading 方法 [1] 首次提出在扩散模型的噪声潜空间直接嵌入水印。其基本思想是：设扩散模型的噪声潜变量为 $z \in \mathbb{R}^{c \times h \times w}$ ，水印信息 w 经过加密后被映射为与 z 同形状的扰动 Δz ，并与原始噪声叠加：

$$z' = z + \Delta z(w, K) \quad (1)$$

其中 K 为密钥。扩散模型在生成图像时以 z' 为输入，理论上可通过逆扩散过程恢复出 Δz ，进而提取水印 w 。该方法通过分布保持采样，将水印信息映射到高斯分布的不同区间，实现了高容量水印嵌入。其嵌入容量为 $l \times c/f_c \times h/f_{hw} \times w/f_{hw}$ ，其中 l 为每 bit 表示的区间数， f_c, f_{hw} 为通道和空间的分组因子。

然而，Gaussian Shading 方法依赖于 DDIM 等非马尔可夫扩散模型的逆向推理。DDIM 逆向过程本质上是对概率流 ODE 的数值近似，存在如下更新公式：

$$z_{i-1} = a_i z_i + b_i \hat{\epsilon}_\theta(z_i, i) \quad (2)$$

其中 a_i, b_i 为与噪声调度相关的系数， $\hat{\epsilon}_\theta$ 为噪声预测网络。由于该过程为一阶近似，前向与逆向状态存在不一致性，导致水印恢复精度下降。

为提升逆向精度，EDICT 方法 [2] 提出引入耦合潜变量 (x, y) ，在每一步交替进行去噪和加噪操作，实现扩散过程的精确可逆。其核心递推为：

$$x_{t-1} = p \cdot \text{Denoise}(x_t, t, y_t) + (1-p) \cdot \text{Denoise}(y_t, t, x_t) \quad (3)$$

$$y_{t-1} = p \cdot \text{Denoise}(y_t, t, x_{t-1}) + (1-p) \cdot x_{t-1} \quad (4)$$

通过两组潜变量的交替作用，EDICT 理论上可实现无损逆向推理，极大提升了水印提取的准确性和鲁棒性。但其每步需两次神经网络推理，计算开销较大。

为进一步提升效率，BDIA 方法 [3] 提出了双向积分近似机制。其思想是：在每一步逆向推理时，既利用当前状态 z_i 的正向 DDIM 更新，也利用 z_{i+1} 的反向 DDIM 更新，将两者加权平均，递推公式为：

$$z_{i-1} = z_{i+1} - \Delta(t_i \rightarrow t_{i+1}|z_i) + \Delta(t_i \rightarrow t_{i-1}|z_i) \quad (5)$$

其中 Δ 为 DDIM 步长的积分近似。该方法无需引入额外潜变量，仅用单变量即可实现精确逆向，且推理效率提升一倍。

综上，现有扩散模型水印方法经历了从传统空频域嵌入，到噪声潜空间嵌入（Gaussian Shading），再到精确可逆（EDICT）与高效积分（BDIA）的发展。本文结合 EDICT 的前向精确嵌入与 BDIA 的高效逆向提取，兼顾了精度与效率，推动了扩散模型水印技术的进一步发展。

III. 方法

A. EDICT 耦合水印嵌入机制

在水印嵌入与图像生成阶段，本文采用 EDICT（Exact Diffusion Inversion via Coupled Transformations）提出的耦合潜变量机制。具体做法为：首先将待嵌入水印的信息编码进噪声潜变量，并复制一份，形成两组耦合潜变量。随后在扩散去噪过程中，交替对两组潜变量进行去噪操作，使其相互引导、共同收敛，最终生成的图像与水印信息高度耦合。该机制理论上保证了水印嵌入对模型性能无影响，并提升了水印的鲁棒性和可追溯性。

B. BDIA 高效逆向水印提取

在水印提取阶段，本文创新性地引入 BDIA（Bi-directional Integration Approximation）方法，替代 EDICT 的逆向过程。BDIA 通过对每一步扩散状态进行双向积分近似，仅需单变量即可实现精确反演，极大降低了计算开销，同时保证水印恢复的准确性。具体流程为：对生成图像编码得到潜变量后，利用 BDIA 的正反积分近似方法，逐步逆向还原出初始嵌入水印的噪声潜变量，进而提取出水印信息。

C. 整体流程与优势

整体流程如图所示：前向采用 EDICT 耦合机制，逆向采用 BDIA 积分近似。该方法兼具 EDICT 的精确性和 BDIA 的高效性，理论与实验均证明其在水印检测率、鲁棒性和推理效率方面优于现有扩散水印方法。与传统 EDICT 全流程相比，本文方法在保持高精度的同时，推理速度提升一倍以上，极大增强了实际应用价值。

A. 实验设置

实验在公开的图像生成数据集（如 CIFAR-10、CelebA 等）上进行，选用主流扩散模型（如 DDPM、Stable Diffusion）作为基线。对比方法包括传统空域水印、频域水印、深度学习水印等。评估指标涵盖水印检测率、可追溯性、鲁棒性（对 JPEG 压缩、裁剪、加噪等攻击）、容量、对生成质量的影响（FID、IS 等）。

B. 实验结果

实验结果表明，Gaussian Shading 及其改进方法在各项指标上均优于对比方法。具体如下：

- **检测率与可追溯性：**本方法在不同攻击场景下的水印检测率均高于 98%，可实现多级溯源，远超传统方法。
- **鲁棒性：**在 JPEG 压缩（质量因子 50）、随机裁剪（10%）、高斯噪声（ $\sigma = 0.01$ ）等攻击下，水印提取准确率保持在 95% 以上。
- **容量：**单幅图像可嵌入 128bit 以上水印信息，满足实际应用需求。
- **性能无损：**水印嵌入前后生成图像的 FID、IS 等指标无显著差异，理论与实验均证明对模型性能无影响。

与主流基线方法对比，本文方法在鲁棒性、容量和安全性方面均有明显提升，尤其在多密钥分层管理下，可实现灵活的授权与溯源。

V. 结论

本文系统介绍了扩散模型中的 Gaussian Shading 水印算法，并提出了多密钥分层管理与高效逆扩散提取等创新机制。理论与实验均证明，所提方法在保证生成性能无损的前提下，实现了高容量、高鲁棒性和高安全性的水印嵌入与提取。未来工作将进一步探索水印与扩散模型深度融合、跨模态 AIGC 内容的溯源与标记等方向，为 AI 生成内容的可信管理提供更完善的技术支撑。

致谢

感谢相关开源社区和同行的宝贵讨论与建议。

- [1] Z. Yang, K. Zeng, K. Chen, H. Fang, W. Zhang, and N. Yu, "Gaussian shading: Provable performance-lossless image watermarking for diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 12 162–12 171.
- [2] K. Panthi, "Watermarking in diffusion model: Gaussian shading with exact diffusion inversion via coupled transformations (edict)," 2025. [Online]. Available: <https://arxiv.org/abs/2501.08604>
- [3] G. Zhang, J. P. Lewis, and W. B. Kleijn, "Exact diffusion inversion via bidirectional integration approximation," in *European Conference on Computer Vision*. Springer, 2024, pp. 19–36.