

[【极简综述】近几年基于深度学习的图像鲁棒水印方法简述 - 知乎 \(zhihu.com\)](#)

[极简综述 - 知乎 \(zhihu.com\)](#)

## 数字图像水印技术综述

### 总起

#### 大概讲了啥

基本模型，特点，分类，攻击以及一些水印算法评价指标

#### 目的

数字水印技术主要用于数字图像的版权保护和完整性认证, 作为一种保护数字图像版权的主要技术, 在数字图像版权受到侵犯时, 能够通过水印提取算法将版权信息提取出来, 作为数字图像归属的主要证据

### 基本模型

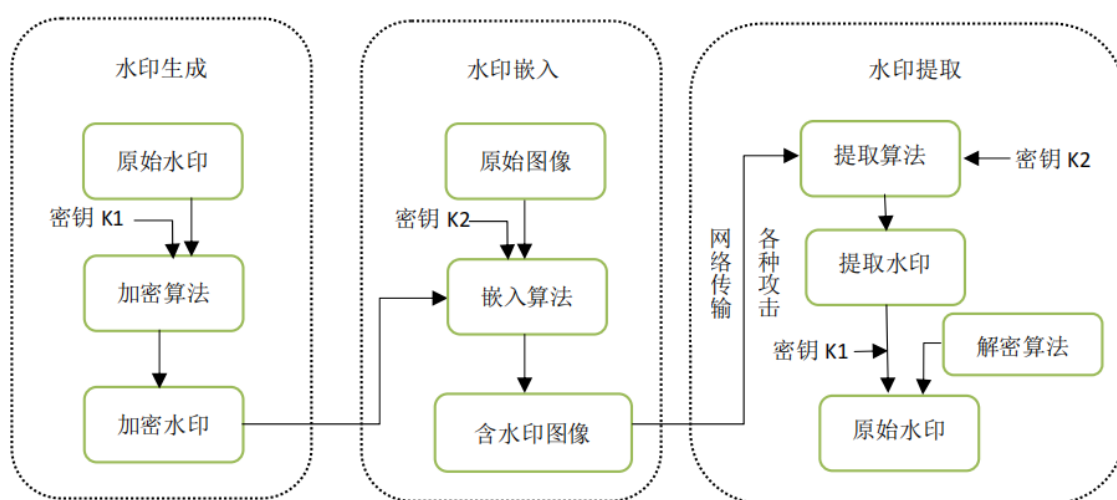


图 1 数字图像水印的基本模型

#### 水印生成

对水印图像加密一般采用置乱方法, 常用的置乱方法有变化模板形状、幻方变换、hash 置乱等

对其他版权信息包括文字信息和数字序列信息加密一般使用伪随机数发生器

#### 水印嵌入

- 空域情况下, 通过密钥 K2 选择嵌入水印的像素位置和像素数量
- 频域情况下, 依据水印信息的数量选择在变换后的高频、中频或者低频域中嵌入水印

#### 水印提取

如果含水印图像在传输过程中被攻击, 提取出来的水印信息就会有不同程度的缺损, 但可以通过计算与原始水印的相似度或者相关系数判断水印信息是否存在

### 数字图像水印特点

- 透明性 (不可察觉)
- 鲁棒性 (水印图像在受到攻击后仍然可以从中提取出水印信息)
- 安全性 (水印能够抵抗恶意攻击的能力)

# 数字图像水印算法分类

## 空域数字图像水印技术

- 最低有效位(LSB)数字水印  
对噪声抵抗力较差, 容易遭受攻击  
改进: 在像素选择方面进行了加密, 可以在奇数行嵌入水印, 也可以在偶数行嵌入水印, 还有根据密钥随机选择像素值嵌入水印
- 二值图像中的数字水印  
通过修改黑白像素个数的奇偶性嵌入水印信息
- 基于图像特征的数字水印  
通过修改原始图像数据使得原始图像的某些统计特征发生变化, 检测时只需要查看含水印图像的统计特征即可, 从而达到盲检测的目的
  - 基于图像亮度值分析的水印算法
  - 基于图像统计特征的水印算法

## 变换域数字图像水印技术

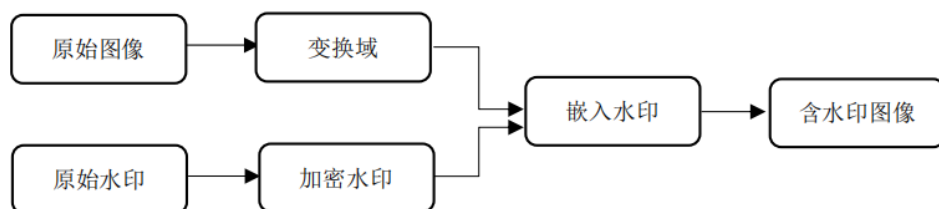


图3 变换域数字图像水印嵌入过程

根据一定的嵌入规则在相应的频带上嵌入水印, 嵌入水印的方式有修改、替换和交换频带的系数

- 低频信息反映了图像的主要轮廓, 不应有较大的失真, 但会**影响到水印的不可见性**
- 高频区域是人眼视觉系统不敏感的部分, 因此在此区域嵌入水印能够保证水印的良好的不可见性, 但是高频区域也是压缩技术常常剔除的部分, 所以**鲁棒性较差**
- 为了满足不可见性和鲁棒性, 一般**将水印嵌入到变换域的中频区域**

主要的水印算法:

- 离散余弦变换(DCT)域
- 离散小波变换(DWT)域
- 离散傅里叶变换(DFT)域

## 数字图像水印算法的攻击

### 鲁棒性攻击

指含水印图像在检测水印之前必须经历的常规信号处理操作

- 非几何攻击 (以削弱原始水印信息强度为主, 对含水印图像进行小幅度的篡改或添加噪声)
  - 噪声攻击: 椒盐噪声、高斯噪声和随机噪声
  - 滤波攻击: 中值滤波、低通滤波和维纳滤波
  - 压缩攻击: JPEG 和 JPEG2000
  - 增强处理攻击: 有锐化、钝化、直方图均衡、Gamma 校正和图像恢复
- 几何攻击 (通过改变含水印图像的像素位置来破坏水印的检测结果)
  - 遮挡、平移和旋转等操作来改变像素的局部或整体位置, 通常包括剪切攻击、旋转攻击、行列偏移攻击和缩放攻击
- 组合攻击 (将几何攻击和非几何攻击进行多种组合的攻击)

## 安全性攻击

指攻击者为了某种利益对水印算法、水印密钥或者含水印图像所进行的各种恶意攻击

## 系统攻击

是针对水印系统中所涉及的其他问题进行攻击

## 数字图像水印技术的评测

### PSNR 峰值信噪比[0,100]

用于衡量嵌入水印后的图像与原始图像之间的失真程度

PSNR 值越大, 说明失真程度越小, 水印的不可见性越好

PSNR 值大于 30 时, 人眼视觉系统不能够感知含水印图像与原始图像之间的差别

$$\text{PSNR} = 10 \lg \frac{\text{MAX}^2}{\text{MSE}}$$
$$\text{MSE} = \frac{1}{mn} \sum_i^m \sum_j^n [f(i, j) - g(i, j)]^2.$$

$f$  表示原始图像,  $g$  表示含水印图像,  $\text{MAX}$  表示图像像素的最大值.

### NC 归一化相关系数[0,1]

用于衡量原始水印信息与被提取水印信息之间的相似程度

NC 值越大, 表示原始水印与提取出来的水印相似度越高, 水印算法的鲁棒性越强

$$\text{NC} = \frac{\sum_{i=1}^m \sum_{j=1}^n (w(i, j) \times w'(i, j))}{\sqrt{\sum_{i=1}^m \sum_{j=1}^n [w(i, j)]^2} \sqrt{\sum_{i=1}^m \sum_{j=1}^n [w'(i, j)]^2}}.$$

$w$  表示原始水印,  $w'$  表示提取出来的水印,  $m$  和  $n$  分别表示水印图像矩阵的行数和列数

### SSIM 结构相似度[0,1]

是对两个图像的亮度、对比度和结构三个量的比较

SSIM 越大, 两个图像之间的相似度越高, 也可以用来衡量压缩图像的质量

亮度度量函数:

$$l(X, Y) = \frac{2\mu_X\mu_Y + c_1}{\mu_X^2 + \mu_Y^2 + c_1}.$$

对比度度量函数:

$$c(X, Y) = \frac{2\sigma_X\sigma_Y + c_2}{\sigma_X^2 + \sigma_Y^2 + c_2}.$$

结构对比函数:

$$s(X, Y) = \frac{\sigma_{XY} + c_3}{\sigma_X\sigma_Y + c_3}.$$

SSIM 函数:

$$\text{SSIM} = l(X, Y)c(X, Y)s(X, Y).$$

其中  $\mu_X, \mu_Y$  分别表示图像  $X$  和  $Y$  的均值,  $\sigma_X, \sigma_Y$  分别表示图像  $X$  和  $Y$  的方差,  $\sigma_{XY}$  表示图像  $X$  和  $Y$  的协方差.  $c_1, c_2, c_3$  为常数, 为了避免分母为 0 的情况, 通常取  $c_1 = (K_1 L)^2, c_2 = (K_2 L)^2, c_3 = \frac{c_2}{2}$ , 一般地  $K_1 = 0.01, K_2 = 0.03, L = 255$ .

# 数字图像鲁棒隐写综述

## 总写

### 大概讲了点啥

基本概念、技术架构、具体方法、性能分析、应用场景与有待研究的问题

### 基本模型

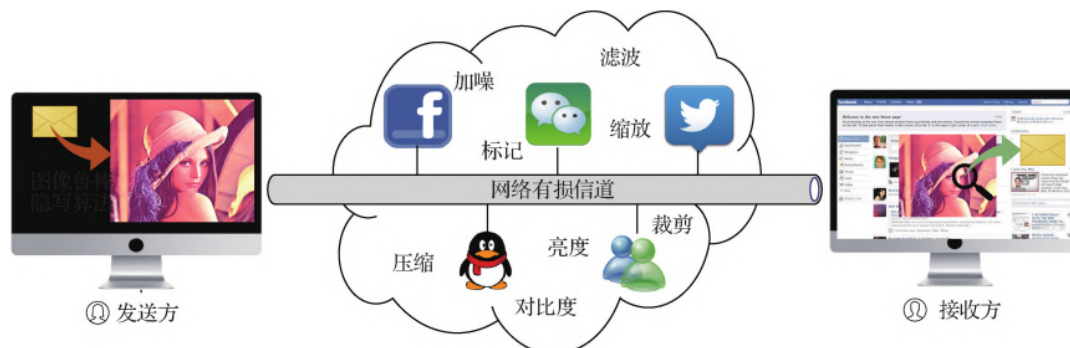


图 1 适应有损信道的图像信息隐藏模型

### 鲁棒水印技术和鲁棒隐写技术的区别

鲁棒水印技术通常关注于水印图像遭受多种信号处理攻击后，水印信息**存在性检测的准确性**，主要通过构造对常见图像处理攻击可保持性较好的信息嵌入域，结合冗余/ 纠错编码等方法，实现低容量、高效率、强鲁棒的水印信息嵌入和提取

鲁棒隐写技术通常关注于载密图像经有损信道传输后，**嵌入信息的不可察觉性和提取准确性**，主要通过设计兼顾鲁棒性与不可见性的信息嵌入方式，结合最小化嵌入代价编码等方法，实现高可靠、大容量、强隐蔽的消息传输

**鲁棒水印技术就是看你水印还在不在，鲁棒隐写技术就是看你水印还能不能提取出来**

## 图像信息隐藏技术

### 图像水印方法

主要关注嵌入水印信息的可恢复性以达到保护版权、证明来源、跟踪盗版的目的

- 脆弱水印（用于数字图像的完整性认证）
- 鲁棒水印（用于数字图像的版权追踪）
  - 基于图像变换的鲁棒水印算法（利用 DCT 变换、小波变换等变换域中系数或关系对常见的图像处理的保持性构造了可抵抗多种攻击的鲁棒信息嵌入域）
  - 基于图像特征的鲁棒水印算法（通过图像特征点/ 区域检测和选择、构建鲁棒信息嵌入域）

### 图像隐写方法

- 经典隐写算法（利用人类视觉对图像某些颜色或系数分量的不敏感性来嵌入信息）
- 自适应隐写算法（选取图像的纹理复杂区域自适应地嵌入秘密信息）
- 可逆隐写技术

## 图像鲁棒隐写架构



## 鲁棒载体构造

- 未知信道边信息

通过**提取**对多种类型及参数的图像处理攻击**鲁棒性较好的隐写载体**、**设计**多种攻击后可保持在相应'0'、'1'有效区间的**载体修改幅度**来实现鲁棒的信息嵌入和提取

- 抗压缩攻击的鲁棒载体构造方法
- 抗几何攻击的鲁棒载体构造方法（主要是针对图像缩放攻击）
- 同时抵抗多种攻击的鲁棒载体构造方法

- 已知信道边信息

- 抗压缩攻击的鲁棒载体构造方法
- 抗几何攻击的鲁棒载体构造方法（主要是针对图像缩放攻击）

## 嵌入代价度量

- 通过对现有经典自适应隐写算法嵌入失真函数的直接引用或简单改进来度量不同鲁棒载体元素的嵌入代价  
仅通过鲁棒载体元素在嵌入信息时的修改幅度反映其在嵌入信息后的鲁棒性与抗检测性能，在为了保持秘密信息在遭受攻击后的可恢复性而对载体元素造成较大更改的情况下度量不准确

- 优化嵌入失真函数并扩展编码方法，改进了 DMAS 算法并提出了 GMAS (generalized dither modulation based adaptive steganography)鲁棒隐写算法

- 通过图像显著性度量优化载体元素嵌入代价

结合社交网络传输图像多包含复杂、显著目标的特点,利用鲁棒图像抽象及显著性检测方法,提出基于复杂、显著区域优先的嵌入代价度量算法

## 嵌入通道选择

通过构造和选择对 JPEG 压缩操作可保持、纹理复杂、不易建模的图像特征区域,实现了鲁棒隐写的嵌入通道选择

基于传输信道匹配的抗压缩隐写算法:

1. 首先辨识信道参数
2. 利用参数识别结果对载体图像进行预处理
3. 通过特定质量因子的JPEG编码器对载体图像进行多次压缩，从而降低信道中的压缩操作对嵌入信息引起的扰动
4. 通过信道匹配、信息嵌入以及重压缩这 3 种操作的反复迭代，直至载密图像在遭受重压缩后仍能完全正确提取其中嵌入的秘密信息

## 信源/信道编码

通常借助最小化嵌入失真编码、差错控制编码等信源/信道编码方法,降低嵌入信息对图像视觉质量的影响,并提高遭受图像处理攻击后信息完全正确提取的概率

## 应用安全策略

- 利用数据分解方法，提出了鲁棒批量隐写方案，使得即使载密图像在有损信道传输的过程中发生了一定程度的丢失，接收方仍可从其余传输完成的载密图像中提取出完整的秘密消息。
- 基于最优载荷分配的 JPEG 图像批量鲁棒隐写方案

## 图像鲁棒隐写方法性能分析

### 鲁棒性测试

对比各种方法在遭受压缩、缩放、加噪、滤波等多种图像处理攻击后,嵌入信息的鲁棒性。

评价指标：平均提取错误率  $R_e$

抗检测性测试

通过 CC-PEV、DCTR 等隐写检测特征,结合集成分类器测试以上算法在不同嵌入比率及编码参数下生成载密图像的抗检测性能。

有待研究的问题

网络有损信道对载密序列影响的精准刻画

兼具多重鲁棒性与不可见性的虚拟载体构建

融合抗攻击、抗检测性的嵌入代价综合度量

高效率、低代价的隐写嵌入通道选择与同步

强容错和低失真共同约束下的隐写编码构造

适应网络有损信道的鲁棒隐写应用策略设计

音频隐写方法综述

概述

基本概念

音频隐写是指利用数字音频或语音作为载体来隐藏秘密信息，主要是在音频信号的时域、频域、小波域以及音频流（Voice of internet protocol, VoIP）中嵌入，关注嵌入信息的隐蔽性与不可感知性，以便保护这些数据免受未经授权者的访问

音频类型

WAV、VoIP、MP3、AMR、AAC、AU、MIDI

表 1 常用音频格式  
Table 1 Common audio formats

类型	格式	特点	适用性
WAV	未压缩	真实记录自然声波形,声音不失真,数据量大	Windows
AU	未压缩	品质音频高、兼容性强、稳定性高	Unix、Java
MIDI	N/A	轻量级、可编辑性、兼容性、音色丰富	多用于音乐制作和演奏等
MP3	有损压缩	对不同频段采用不同的压缩率,压缩后占用空间小	多平台适用,常用移动设备
AAC	有损压缩	编解码器质量高,性能高	多平台适用
AMR	有损压缩	压缩比较大,满足移动通讯需求,对于通话效果较好	多用于移动设备
VoIP	有损压缩	适用网络通讯,减少失真	智能手机与电脑在内的联网接入设备

工具

表 2 音频隐写工具  
Table 2 Audio steganography tools

名称	支持格式	算法	开源
MP3Stego	MP3	基于量化步 长修改方法	✓
S-Tools	WAV	LSB	✓
SilentEye	WAV	LSB	✓
Hide4PGP	MP3/VOC	LSB	✓
stegandomet	MP3/WAV/MIDI/AU	LSB	✓
StegoStick	WAV	LSB	✓
Info Stego	MP3	N/A	✓
Scram Disk	WAV	N/A	✓
DeepSound	MP3/WAV	LSB	✓
Steghide	WAV/AU	LSB	✓



评价指标

- 不可感知性（以下指标均为越高越好）
  - 信噪比 SNR （表示含密音频的失真度）
  - 峰值信噪比 PSNR （表示含密音频的失真度）
  - 客观等级差异 ODG
  - 均方误差 MSE （表示含密音频的失真度）
- 隐写容量（在满足不可检测性的条件下分析隐写容量的极限）
- 鲁棒性
  - 比特误码率 BER （错误比特数与总比特数之比）
  - 等错误率 EER （衡量是否保留了说话人的身份信息，EER越低身份信息越完整）

传统音频隐写方法

主要分类

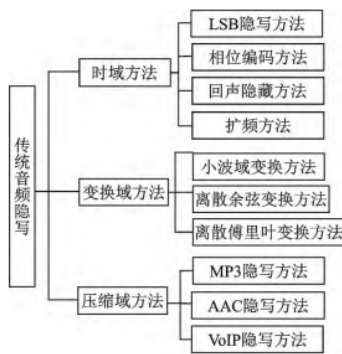


图2 传统音频隐写方法  
Fig.2 Traditional audio steganography methods

对比

表4 传统音频隐写方法对比			
Table 4 Comparison of traditional audio steganography methods			
嵌入域	方法	优点	缺点
时域	LSB隐写方法	算法简单,易于实现、计算复杂度低、隐藏容量大	鲁棒性差,抗检测能力弱
	相位编码隐写方法	鲁棒性好、可以有效调节不可感知性	具有一定的相似度
	回声隐写方法	算法简单、隐藏效果好、不产生噪声、能够实现盲检测	隐藏容量小、提取效果差、信道噪声影响大
	扩频隐写方法	鲁棒性好、不可感知性好	隐藏容量小、算法相对复杂
变换域	FFT隐写方法	稳健性好、隐写容量大	不可感知性差
	DCT隐写方法	隐写容量大	计算复杂度高
	DWT隐写方法	隐写容量大、不可感知性好	计算复杂度高
压缩域	MP3隐写方法	不可感知性好	隐写容量小
	AAC隐写方法	抗隐写分析较好、鲁棒性好	高比特率下的隐写容量小
	VoIP隐写方法	实时性好、隐写方法灵活、隐写区域多、抗检测性好	隐写容量小

基于深度学习的音频隐写

嵌入式载体式音频隐写（在数字音频上通过深度学习方法完成秘密信息的嵌入和提取）

- Encode-Decoder 结构
    - 基于GCNN的音频隐写方法
- $C$  为载体音频,  $M$  为秘密信息,  $E$  表示载体音频编码器,  $H = [E(C), C, M]$  为三者的连接, 通过含密载体解码器  $D_c(\cdot)$  得到含密载体频谱  $\tilde{C} = D_c(H)$ , 再利用  $c$  的相位  $\angle c$  通过  $S'(\tilde{C}, \angle C)$  得到语音,  $S(\cdot)$  为 STFT. 接收方将接收到的语音经过  $\tilde{C} = S(S'(\tilde{C}, \angle C^*))$  得到含密载体  $\tilde{C}$ ,  $S'()$  为逆 STFT, 最后通过  $D_m(\cdot)$  秘密信息解码器获取重建秘密信息  $\hat{M} = D_m(S)$



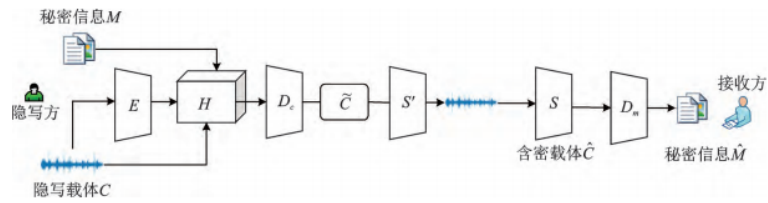


图6 基于GCNN的音频隐写方法(Hide & Speak)

损失函数:  $L(C, M) = \lambda_c |C - \hat{C}| + \lambda_m |M - \hat{M}|$

缺陷: 难以在一个混合音频中的单个音频中隐藏信息

- 基于DNN的源混合和分离的音频隐写模型 (MSRAS模型)

接收方通过源分离的操作得到  $\tilde{c}_i$ ，再通过解码器重建秘密信息

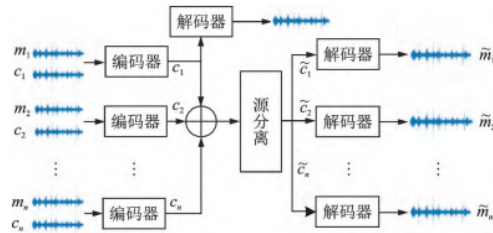


图7 MSRAS模型

- 可以在音频中嵌入图片信息的音频隐写方法 DITAS

利用多级网络将秘密图像的残差分阶段逐步嵌入到多个音频载体中，利用基于 U-Net 省略全连接的 2D 全卷积神经网络构成编码器和解码器，编码器对需要隐藏的图像进行编码，将其添加到音频载体 STFT 的频域中

缺陷: STFT (短时傅里叶变换) 同时具有幅度和相位变换, 存在相位重建问题

- 基于 STDCT (短时离散余弦变换) 的新型残差网络结构的音频隐写模型 PixInWav

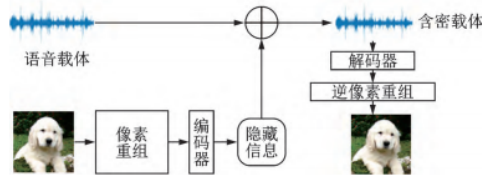


图8 PixInWav 模型

- 自动学习嵌入代价

- 利用 GAN 实现嵌入

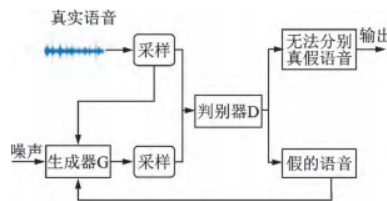


图9 GAN模型

- 基于深度卷积 GAN 的音频隐写模型 DCGAN

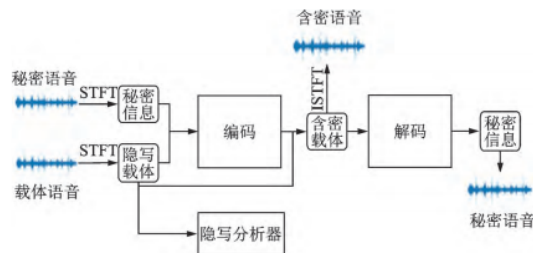


图10 基于深度卷积GAN的音频隐写

- 批处理归一化优化 SNGAN 的隐写方法

针对现状: 基于 GAN 的音频隐写方法忽略了隐写容量和不可感知性的高要求

在 GAN 的生成器和判别器中分别利用了频谱归一化，编码器、解码器和隐写分析器三部分协同学习

o 基于 MIDI 和 GAN 的音频隐写方法

针对现状：音频隐写的不可感知性和抗检测性差

利用 Music21 工具包构建带有索引的 MIDI 音符字典，利用 GAN 网络的生成器、提取器和判别器网络进行训练

• 基于对抗样本的方法

？：这里的嵌入方法 LSBM 是 LSB 吗（

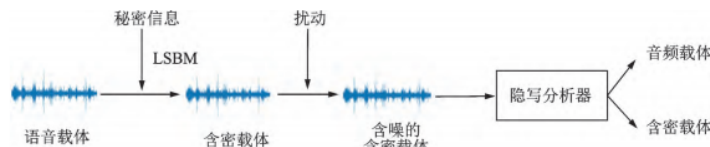


图 11 基于对抗样本的音频隐写

对比：

表 5 嵌入载体式音频隐写模型对比

Table 5 Comparison of embedded cover-based audio steganography models

模型	嵌入域	时频变换	秘密类型	模型特点	评价指标	数据集
Hide & Speak <sup>[67]</sup>	频域	STFT	音频	基于门控卷积神经网络的音频隐写模型，不可感知性好，且可以在单个音频中隐藏多个秘密信息	SNR	TIMIT YOHO
DCGAN <sup>[80]</sup>	频域	STFT	音频	基于 GAN 的音频隐写模型，具有少量参数的轻量级模型，可用于物联网的许多设备中	SNR	TIMIT Librispeech
U-NetGAN <sup>[79]</sup>	时域		音频	利用 GAN 实现音频隐写在时域内的最优嵌入，抗隐写分析能力强	Error rate	UME-ERJ WSJ
CNNAE <sup>[86]</sup>	时域		音频	基于对抗样本的音频隐写模型，不依赖于现有的隐写成本，抗隐写分析能力强	Accuracy	TIMIT
BNSNGAN <sup>[81]</sup>	时域		音频	可以实现任意长度秘密音频的嵌入，并且在不可感知性、隐写容量和抗检测性上有较好的均衡	SNR ODG BER	TIMIT Librispeech
MSRAS <sup>[71]</sup>	时域		音频	对源混合和分离具有鲁棒性的音频隐写，秘密信息被单独隐藏在某些源中，并再与其他源混合和源分离后能够准确恢复秘密信息，鲁棒性强	SNR	MUSDB18
DITAS <sup>[72]</sup>	频域	STFT	图片	利用多级网络将图片隐藏到音频中的多模态隐写方法，有效载荷容量控制灵活，隐藏容易	MSE PSNR	TIMIT LJ Speech VOC2012
PixInWav <sup>[74]</sup>	频域	STDCT	图片	基于 STDCT 的多模态隐写方法，能够在不影响隐写载体质量的情况下独立编码图像，并可以离线编码图像	SNR SSIM	FSDnoisy18K ILSVRC2012
MIDI-GAN <sup>[84]</sup>	时域		音频	突破有载体隐写在不可感知性和抗隐写检测性的限制，将秘密信息转化为 MIDI 音频，从而提高载密音频的有效性安全性	MOS	MIDI
PixInWav2 <sup>[77]</sup>	频域	STFT	图片	对 PixInWav 模型损失函数的修改，STDCT 替换为 STFT，以及在编码过程中引入冗余进行纠错等增强鲁棒性	SSIM PSNR SNR	ILSVRC2012 FSDnoisy18k
LSBMAE <sup>[88]</sup>	时域		音频	在 LSBM 方法嵌入得到的含密音频载体上加入扰动并通过训练好的隐写分析仪进行误分类，具有高感知的含密载体，并且抗检测性强	PSNR Accuracy	TIMIT UME

## 生成载体式音频隐写

利用深度神经网络生成适合隐写的音频载体，然后在生成的载体上完成秘密信息的嵌入和提取

训练判别器来逼近生成的语音载体和真实语音的相似性

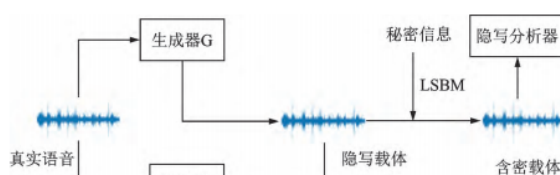


图 12 基于 GAN 的生成载体式音频隐写

无载体式音频隐写

由隐写模型根据秘密信息直接生成含密载体

- 基于 RNN 的无载体式音频隐写

在音频生成过程中，可以根据每个音符的条件概率分布对其进行合理编码，然后根据位流控制音频生成。同时，可以通过精细调节编码部分来控制信息嵌入率，从而实现隐蔽性和隐藏容量的同步优化。

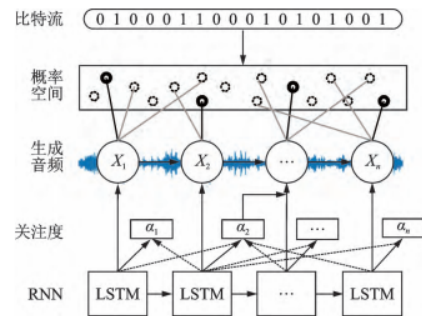


图 13 基于 RNN 的无载体音频隐写

- 基于GAN 的无载体音频隐写方法

利用音频合成模型 Wave GAN 作为生成模块的基础，并将输入的秘密音频直接生成为含密载体音频实现生成式隐写

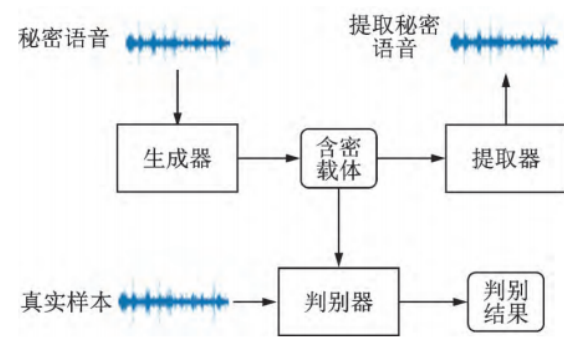


图 14 基于 GAN 的无载体音频隐写

Fig.14 Coverless audio steganography based on GAN

分析对比

表 6 基于深度学习的音频隐写方法对比			
Table 6 Comparison of audio steganography methods based on deep learning			
方法	实现过程	优点	缺点
利用 LSTM 和 Attention、GAN			
生成载体式音频隐写	生成载体音频,再通过哈夫曼树的嵌入与提取	无需准备音频载体,能够直接生成音频载体	受模型及训练过程的影响,生成音频载体质量不高
嵌入载体式音频隐写	利用 CNN、GAN 等自动学习实现秘密信息的嵌入与提取	自动学习嵌入代价、秘密信息的嵌入与提取	计算效率低,秘密音频的损伤
无载体式音频隐写	利用 GAN 模型根据秘密信息直接生成含密音频	能够根据秘密音频直接生成含密载体,抗检测能力强	受秘密音频的限制,隐写容量低

进一步研究方向

基于深度学习的音频隐写方法

鲁棒音频隐写方法

轻量级隐写方法 (计算复杂度更小)

行为安全的音频隐写方法

语音实时隐蔽通信系统