

文档隐写溯源技术分析

原创

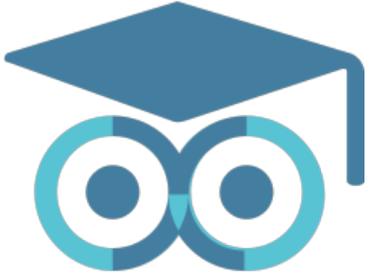
拾年之璐 于 2021-09-25 23:08:39 发布 921 收藏 7

分类专栏: [研究生课程](#) 文章标签: [隐写](#) [溯源](#) [文档隐写](#) [隐写溯源](#)

版权声明: 本文为博主原创文章, 遵循 [CC 4.0 BY-SA](#) 版权协议, 转载请附上原文出处链接和本声明。

本文链接: https://blog.csdn.net/cxh_1231/article/details/120479482

版权



[研究生课程](#) 专栏收录该内容

19 篇文章 16 订阅

订阅专栏

本文是信息安全综合设计的课程报告。全文字数: 4000+。



本文主要内容:

- 1 文档隐写溯源技术的意义
- 2 当前文档隐写溯源技术存在的不足
 - 3.1 文档管理系统的改进
 - 3.2 文档隐写技术的改进
- 4 改进文档隐写溯源技术效果分析

1 文档隐写溯源技术的意义

在如今的信息时代, 信息技术的迅速发展使原本以纸张形式保存的大量重要文档转变为以电子文档形式保存。得益于电子文档的创建快速、易于存储、方便管理、传播共享等优点, 许多企业、政府单位的日常文件甚至机密信息都是利用电子文档的形式进行存储与传输。这些文档中包含了多种多样的信息, 具有巨大的经济价值和应用价值。

但是，许多涉密文档的使用，并没有得到有效的监督和管理，信息泄露事件频有发生。2021年5月9日，宁夏银川一份名为《关于境外输入无症状感染者XXX的初步流行病学调查报告》的流调报告遭外泄，其本人及家人详细信息被曝光。2021年1月14日，一份名为《河北藁城确诊人员基本情况》的文件截图，在微信群内迅速传播，其包含全部人员的详细信息；2020年12月23日，辽宁沈阳一新冠肺炎阳性病例个人信息遭泄露；2020年12月7日，四川成都新冠肺炎确诊病例患者赵XX的流调报告遭外泄，泄露者于9日被行政处罚……从2020年11月以来，全国各地至少出现10起流调信息泄露事件，泄露之后，虽然各地政府介入调查，但是只有极少数事件能追溯到泄露的源头，大部分事件的溯源工作都无结而终。以上只是近期发生的，在网络上产生广泛舆论的文档泄露事件。除此之外，还有更多的涉密文档泄露事件发生。

观察这些涉密文档泄密事件，可以看出泄露的途径，主要是被文档的合法使用者非法使用，通过源文件转发或者截图转发的形式，将涉密信息发布于公众平台，从而导致广泛传播。这种源文件不具有唯一性和可溯源性，在传播过程中，无法得到有效的监督和控制。

针对上述缺陷，拟采用一种文档隐写溯源技术，对涉密文档进行处理后传输，保证每份文档的唯一性和扩散路径可控可溯源，能够高效快速查询文档的使用和扩散过程，防止文档被滥用。当泄密事件发生时，能够快速定位溯源，为追责提供依据，从而提高涉密文件的安全性。另外，文档隐写溯源技术的采用，在版权保护等方面，也有重要的研究意义。

2 当前文档隐写溯源技术存在的不足

在当前的文档溯源技术中，最常用的是采用明文水印（特定标识）的方式。一种是使用统一的水印，但这种形式的水印，有与无并无差异，无法起到溯源的作用。另一种是水印差异化，即不同用户打开同一个文档时，文档显示的水印通常是用户的ID或者用户名。不可否认，这是一种简单的防止涉密文档外泄的有效手段。比如常用的OA办公系统钉钉，可以开启群文件的“保密模式”，即只支持在线预览，并带有水印，但不可下载就限制了文档的打印。此外，这种明文水印方式，使得水印可视，对文档的阅读效果有影响，如果采用技术手段将水印去除，则泄密后，也无从溯源。

数字水印技术是一种版权保护技术，其主要应用于图像、视频等载体，后期也用于文档中。该技术通过对文档的某种特定结构（如行间距、字符间距等）进行修改，用来存储信息，从而实现水印信息的嵌入。嵌入后的文档包含了水印信息，但是这种数字水印信息是不易被察觉和修改的，并且文档的原价值不受任何影响。采用这种技术嵌入的水印，只能被嵌入者识别和提取，能够起到文档隐写和溯源的功能。但是这种数字水印技术加密后的文档，一方面是存储信息的容量低，另一方面是鲁棒性和可用性差，即如果对文档图像进行了缩放或者翻转，字符间距等信息可能改变，使得字符间距达不到可识别的阈值。

针对当前文档隐写溯源技术中存在的不足，可以从文档管理系统和文档隐写技术两方面进行改进。

3.1 文档管理系统的改进

当前常用的OA办公系统，并没有对文档采用加密溯源等技术，只是采用简单的明文水印限制扩散，但是也限制了文档的再次编辑。所以，可以将当前的文档管理系统分成四个模块来改进加密溯源技术，分别是文档上传模块、权限控制模块、文档控制追踪模块和记录查询模块。这几个模块的详细功能如下。

（1）文档上传模块

该模块的主要功能是进行文档的上传，这也是原始文档的最初来源。

（2）权限控制模块

该模块主要是生成上传的文档的访问权限列表，文档发布者可以设置哪些用户可以对文档进行访问，并且访问的权限又分为不可访问、只读、只读可复制、可批注、可编辑以、完全访问六种权限。

（3）文档控制追踪模块

该模块是文档溯源技术改进的核心，主要实现文档的访问控制，一方面根据访问权限列表，生成指定权限的文档，分发给用户，另一方面记录用户对文档访问痕迹的溯源记录，以及记录文档扩散的痕迹，生成完整的溯源链。这里的溯源记录信息，通常包含此条访问记录的标识符、访问此文档的用户信息（用户名）、访问时间、访问硬件设备地址、文档名、该溯源记录的前置溯源标识符等信息。除了访问父文档会产生溯源记录，在此OA系统中发送过的任何子文档，都会产生溯源记录，所以这里的前置溯源标识符，就是用来记录上一个文件的溯源记录标识符。生成的溯源记录，一方面将详细信息存储于文档管理系统服务器，另一方面将溯源信息以不可见的形式，嵌入到分发给用户的文档中。

此外，该模块还应该实现加密和签名功能。通常采用公钥加密、私钥解密和私钥数字签名、公钥验证等机制保证溯源记录的保密性和真实性。再此过程中，为保证文档访问时间的准确性，还需要增加数字时间戳服务中心，来对数字时间戳进行验证盖章。

(4) 溯源查询模块

该模块主要进行溯源记录的查询。

一方面可以查询指定文档的扩散链。如Alice上传了文档D1，Bob下载了文档D1，编辑为D2后，通过OA发送给了Carol。这时可以查询文档D2的扩散链：Alice→Bob→Carol。

另一方面可以查询任意文档的扩散链。如Carol将Bob通过OA发来的文档D2发到了外网，这时，Carol下载文档已变成D3。Alice发现了外网的文档D3，Alice可以通过上传D3文档到查询模块查询该文档的详细扩散链，其中系统解析文档溯源记录标识符，进行匹配查询。

3.2 文档隐写技术的改进

通过对文档管理系统的改进，可以实现对任意文档的溯源。但是，如果Carol通过截图，将部分文档内容发送至外网，此文档管理系统的溯源功能就失效了。除了使用基本的明文水印，还可以通过对文档进行隐写，即改进隐写技术实现溯源功能。改进的隐写技术可以从两方面入手：盲水印或文字欧拉数。

(1) 盲水印

一般来说，上级下发的涉密文档是不可编辑的，所以可以通过将文档转成图片形式，并添加盲水印的方式，实现文档隐写。通过这种方式添加的水印，是肉眼不可见的，隐匿性和抵抗攻击能力强。

目前常用的添加盲水印的方式可分为空域方法和频域方法，在文档图像中添加冗余信息（水印），并保证图像质量不变。常见的图像是空域的。空域方法是在空间域直接对图像进行操作，水印叠加在图像上。而频域方法是通过某种变换手段，比如傅里叶变换、小波变换等，将图像变到频域，再对图像添加水印，在进行逆变换，将图像转换成空域，再保存带有水印的图像。图一、图二为频域方式加水印和提取水印的基本流程。

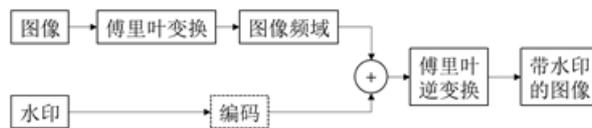


图1 频域方式添加盲水印基本流程



图2 频域方式提取水印基本流程

相比较于空域，频域方式添加盲水印的隐匿性更强，抗攻击性更高。下图是采用傅里叶变换实现的添加盲水印的示例。



图3 添加与提取水印示例

(2) 文字欧拉数

采用添加盲水印的方式对文档进行加密，可以有效对通过截图方式泄密的文档溯源。而对于摄屏、打印等方式泄密的文档，盲水印的方式就失效了。为此，可以通过对文字的欧拉数进行变换，来实现文档隐写。

欧拉数（Euler Number）是一个工程中常见的参数，其具体意义在不同的学科中是不同的。在拓扑学中，欧拉数通常表示空间的完整性。对于一个文字来说，其欧拉数可以采用拓扑学中的欧拉函数来计算。其计算方法如公式(1.1)所示。

$$EUL = C - H \quad (1.1)$$

式中，EUL表示欧拉数，C表示对象数（连通区域数），H表示空洞数量。

比如图4（左图）中的“理”字，其“王”字旁和右侧的“里”是分开，不相连的，故其对象数为2，右侧的“里”有四个空洞，故其空洞数量为4，所以“理”字的欧拉数EUL为-2。



图4 左：黑体“理”字，右：改变拓扑结构的“理”字

对于同一种字体，同一个汉字的欧拉数是固定的。那如何通过欧拉数实现信息的隐写呢？可以通过改变汉字的笔画，来改变汉字的欧拉数。

再如图4（左图）的“理”字，将“王”字旁的第一横和竖之间做拆分，那么“王”字旁的对象数就是2，新的“理”字的对象数为3，欧拉数为-1。这就实现了欧拉数的改变。

除了可以计算单个汉字的欧拉数，还可以根据一串汉字的欧拉数，来实现信息的隐写。如图5中的四个“武汉理工大学”字符串，乍一看，没有什么不同。仔细一看，①和②的“武”字不同，改变了“武”字的对象数；②和③的“汉”字不同，改变了“汉”字的空洞数；③和④的“学”字不同，改变了“学”字的对象数。

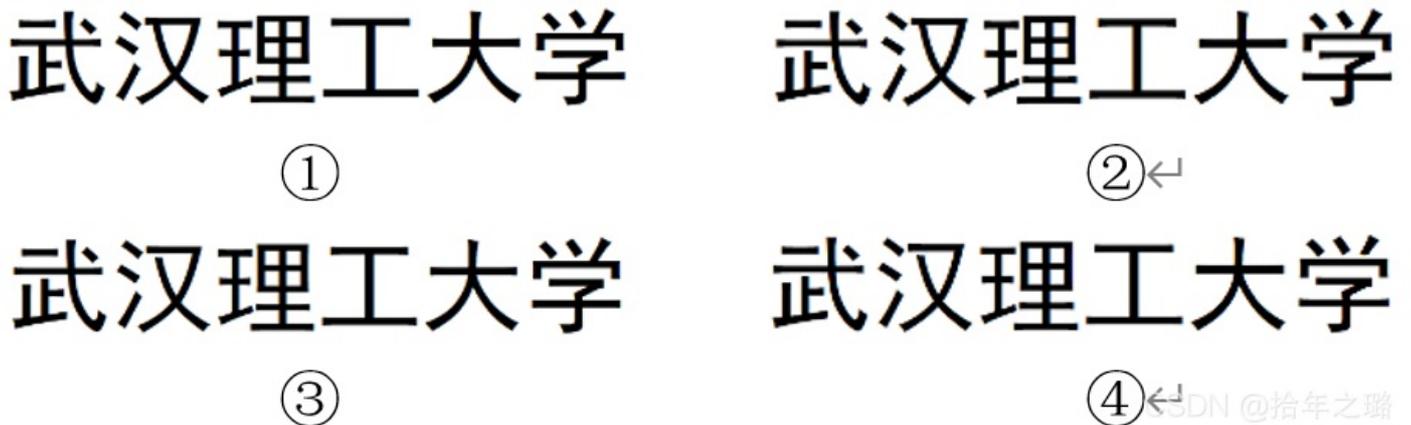


图5 四个不同欧拉数的“武汉理工大学”

使用MATLAB可以轻松计算出这四个“武汉理工大学”字符串的欧拉数，如图6所示。



图6 计算字符串的欧拉数

有了文字的欧拉数，可以使用该技术来改进文档隐写技术。具体的操作过程如图7所示。对于每一次请求获取的文档，首先使用图像识别算法，识别出各个文字的欧拉数。然后此处选择每一行的前6个文字（非标点符号）的欧拉数组成行向量。如果此行的序列在数据库中已存在，则随机选择几个汉字，改变其拓扑结构，如图7的红线所示。改变拓扑结构的技术可以使用二值图像技术。最终将每行的行向量组合一个n行6列的矩阵，矩阵的每一行的行向量都是不同的。对于任意一份文档，或者是文档中的部分截图，可以通过识别截图中的几行汉字的欧拉数，组成矩阵，与数据库中的矩阵进行对比溯源。

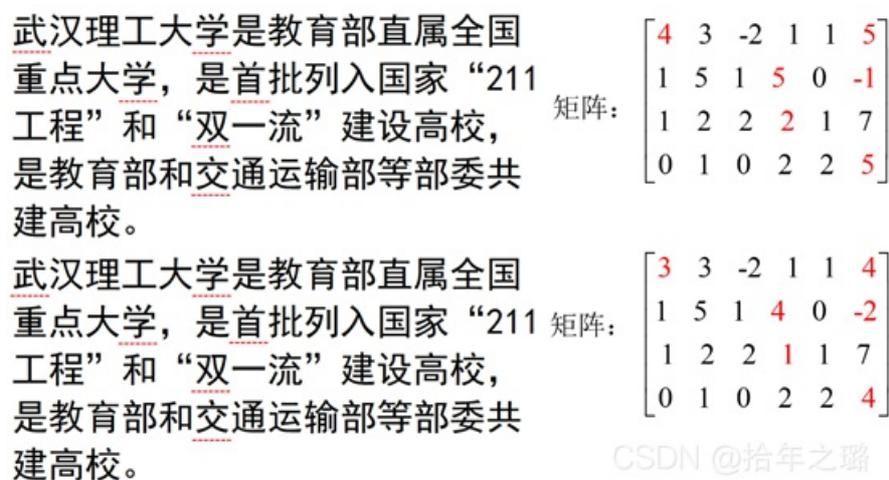


图7 根据欧拉数计算文本矩阵（取前4行，每行取前6个汉字）

4 改进文档隐写溯源技术效果分析

本文通过两个方面的改进，对文档隐写溯源技术进行改进。每个方面都有自己特点，也有自己的针对点。

对文档管理系统的改进是针对文档本身作保护。这不仅能够严格控制文档的扩散过程，显示完整的扩散链，而且能够快速定位泄露文档的源头，保证每一份文档都能够得到有效的监督和管理。

通过对文档隐写技术的改进，是保证文档的截图、照片、打印件可溯源。使用欧拉数来记录每一份文档、差异化每一份文档，当文档被截图泄露时，能够通过提取文档中的水印，根据水印信息来确定文档的来源，完成泄密文档的溯源。

总之，通过此文档隐写溯源技术的改进，提高了涉密文档的安全性，无论是在文档本身还是在文档的截图、纸质方面。此项技术的应用，能够有效震慑接触涉密文档的人员，对于泄密文档的溯源追责有强有力的帮助。