

# 基于 Merkle 树的移动平台文件完整性校验<sup>①</sup>

张晓燕 范冰冰 (华南师范大学 计算机学院 广东 广州 510631)

**摘 要:** 移动平台的下载业务是除消息类业务外另一类重要的移动增值数据服务, 对下载过程的有效管理和控制是该业务稳定运营的前提。下载数据的完整性是下载业务的核心内容。描述了移动平台下载业务内容完整性校验方案。提出了 Merkle 树在移动平台内容完整性校验上的实现策略, 该策略通过缩短校验值在网络传输中的大小, 来有效减少完整性校验的代价。理论分析与工程项目证明, 该机制具有较高的安全连通性和较低的网络延迟。

**关键词:** 完整性校验; 移动平台下载业务; Merkle 树

## File Integrity Check on Mobile Platform Based on Merkle Tree

ZHANG Xiao-Yan, FAN Bing-Bing (Dept. of Computer, South China Normal University, Guangzhou 510631, China)

**Abstract:** In addition to the message service, the download service is another important business. The effective management and control of the download process is a prerequisite for stable operations for the business. The data integrity is the core content for download business. This article describes the content integrity check program in the download process of a mobile platform. It proposes a content integrity check strategy on the mobile platform based on Merkle tree, which reduces the cost of the integrity check effectively by shortening the size of checksum on the course of the network transmission. Theoretical analysis and the project has proven that the mechanism has high security connectivity and low network latency.

**Keywords:** integrity check; download business on mobile platform; Merkle tree

## 1 引文

随着 3G 网络的飞速发展, 移动平台应用出现多样化, 多元化的发展。移动数据业务已经成为研究和实现的热点。移动增值业务、移动互联网、多媒体短消息等, 被认为是移动通信产业新的利润增长点。随着移动数据应用的不断丰富, 内容下载类业务也正在得到越来越多的应用。用户可以从相应内容服务器端下载不同的铃声、屏保、音视频片段用于本地播放或欣赏, 升级或安装各种客户端软件用于本地或联网应用等。由大唐电信推出的基于短消息(SMS)的 GSM 手机动态 STK 业务空中下载技术, 动态 STK (SIM Application Toolkit), 是在 GSM 11.14 中提出的一

种开发工具。STK 采用基于短消息的机制, 实现了部分的数据业务由 PC 转到手机, 满足了用户在移动中获取信息的需要。动态 STK 业务空中下载技术采用了先进的 OTA (Over the Air) 技术。使用者在 GSM 网络覆盖的范围内可以随时随地下载新业务。业务空中下载技术解决当前 2G 移动通信网络增值业务更新的最佳方案之一。但是, 该技术是基于短消息的下载技术, 下载的信息量严重受到限制。

笔者参与的项目《移动商店平台》移动应用下载业务的研究(如图 1 所示), 出发点是解决目前移动应用下载业务的应用瓶颈, 注重集中在获取下载资源、传输下载资源、下载业务管理控制这个大的方面, 获

<sup>①</sup> 基金项目: 广州科技支撑计划(200922-D261)

收稿时间: 2010-01-04; 收到修改稿时间: 2010-03-24

取下载资源。在该项目中,无线应用下载业务作为一个独立的模块,实现了基于动漫、小说、新闻、音乐、图片等内容的下载应用。完整的无线应用下载业务应该能使无线用户在网上找到感兴趣的内容,通过从无线终端,如手机发起下载请求,实现为完成该下载过程所需要的基础功能,如多线程下载、断点续传等功能,更为重要的是提供一些机制保证这些基础功能得以实现,如内容完整性校验。为解决移动业务应用下载中数据块校验的问题,在移动网络中,应用下载有其特殊场景,大多为集中控制与集中下载,同时所请求资源也相对集中。在此情景下,文件分块下载是常用手段,实现文件完整性校验以尽早发现与定位出现问题的块,从而减少重传的资源块是提高下载性能最重要的手段之一,在目前的研究背景下,笔者提出一种优化策略,以解决现有技术中应用面临的窘迫。

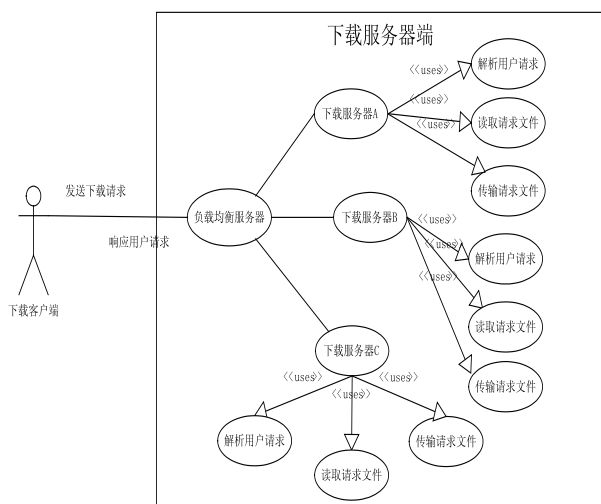


图1 移动应用下载业务用例图

文件完整性检测的根本思想是:将被入侵的系统状态和未被入侵的系统状态相比较。在建立了正确的初始化系统状态后,检测程序定期地根据初始系统状态检测当前系统状态,发现可疑或非法的变动,触发警报和通知系统管理员。实现一个完全的检测功能,至少需要几个部分:正确安全的系统状态、预先定义的被检测系统状态和一个进行定期检测的程序<sup>[1]</sup>。具体到文件检测时,现有的方法是采取密码学的方法检验文件的完整性。文献[2]阐述:“数字签名,是最直接的一种内容完整性校验方法。数字签名是数字证书的基本应用之一,数字签

名的基本过程是信息发送者先对原始数据进行杂凑运算(即 HASH 运算)得到消息摘要,再采用自己的签名私钥对消息摘要进行加密运算。验证签名的基本过程是信息接收者对收到的原始数据采用相同的杂凑运算得到消息摘要,并使信息发送者的签名公钥对数字签名进行解密获得消息摘要,将二者进行对比,以校验原始数据是否被篡改。”数字签名可以完成对数据完整性的保护和传送数据行为不可否认性的保护。相对比杂凑函数而言,使用 RSA 算法签名验证的计算量非常大,由此不适合实时性要求比较高的应用。

文献[3]提出:目前,在互联网应用领域,采用数字签名校验数据完整性的一般过程。其中,签名私钥配合散列算法的使用,可以完成数据签名的功能。在数字签名过程中,可以明确数据完整性在传递过程中是否遭受破坏和数据发送行为是否是签名证书所声明的身份的行为,提供数据完整性和行为不可否认性功能。

文献[4]提到“文件完整性检验方式:根据用户定制的配置文件中需要校验的文件系统内容进行散列计算,将生产的散列值与文件完整性数据库中存储的预先计算好的文件内容的散列值进行比较。不一致则说明文件被非法更改,并可判定发生入侵。”常用的散列算法有 MD5, CRC16, CRC32, MD2, SHA。

文献[5]分析了 MD5 用于文件完整性检测的优点,并通过实验证明基于 MD5 的文件完整性检测软件是高效、实用的。JAVA 提供了 MD5 的实现。笔者在项目中采用 MD5 算法来计算文件的信息摘要。在移动应用平台上,下载业务需求:分块下载,断点续传,对每个分块都要进行完整性校验。这就必然要求,一个文件有 1~N 个文件块(N 需要根据具体的传输最大字节数和文件大小确定),每个文件块在传输的时候,除了自身数据,还要传输数字签名。如果有一种方式,能够解决:

(1) 对服务器而言,一个完整的文件下载过程,服务器只需向客户端发送一次数字签名,就能够校验文件的完整性;

(2) 对客户端而言,没有必要在校验文件中保存所有数据块的校验值了,仅需保存其中计算需要的哈希值即可。

这样做的好处很明显,减少网络的负载量,增强

网络传输的实时性。对移动应用业务而言,这些改进很关键。

## 2 基于Merkle树实现文件完整性校验

文献[6]阐述了:1979年由Merkle提出的一种树结构:Merkle树=二叉树+HASH算法,merkle树是一种特殊的二叉树,其集合中的元素作为树的叶子结点,树的每一个内部结点是其左右孩子级联后的HASH值,最后再以签名等方式认证根节点。

对于一个Merkle树,可以构造一条从叶子结点到树根节点的HASH链表,HASH链表的每一单元中都包含了一个元素和一个位置标识 $o_i$ ,该标识用于表示元素应该从左侧(L)还是右侧(R)进行级联。

Merkel HASH树是一种高效的用于数据验证的索引结构。Merkle HASH树与简单的数字签名方法不同的地方在于,简单的数字签名方法对每个元组都计算一个数字签名,而基于Merkle HASH树的方法在分块数据之上构建一个二叉树,二叉树的每个叶子节点对应一个数据元组,每个叶子节点存储一份对应数据元组经过一个单向HASH算法计算机以后的HASH值,对应的,每个中间节点存储一份由其孩子节点HASH值计算机以后得到的HASH值。最后,对二叉树根节点计算一个数字签名,用于之后验证数据的正确性以及完整性。

### 2.1 基于Merkle树实现文件完整性校验的过程

依据Merkle树定义,构造用于验证移动下载业务文件完整性校验的Merkle树,其中HASH算法使用MD5,如图2所示。

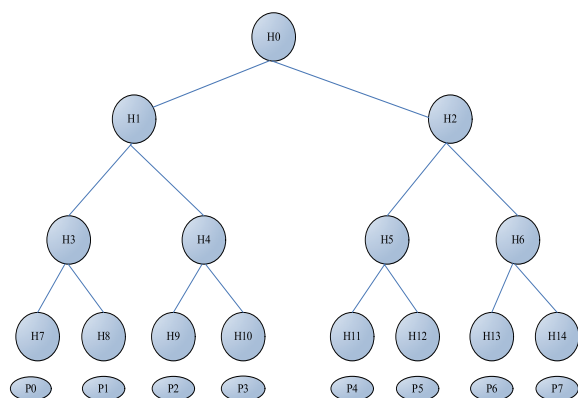


图2 Merkle MD5树

假设某文件共7块数据P0~P6,根据Merkle树

需要将其填充到8块,填充的P7仅用于填充校验,而无需当作数据块传送,每个块均有对应一个HASH值,对于每个父节点,将两个子节点的哈希值级联后求出HASH值作为父节点HASH值,依此类推,直到求出根节点的MD5值,H0,这已计算过程就构成了一颗二元的Merkle哈希树,树中的叶子节点H7~H14对应着数据块P0~P6的实际哈希值,非根结点和叶子节点的其他所以节点称为路径哈希值,即H1~H6。

如图3所示,当用户提交一个下载文件请求,服务器除了返回用户对应的响应数据外,还需要返回验证对象相关信息和Merkle HASH树根节点的数字签名。用户结合接受到的响应数据和验证对象,在本地计算机出Merkle HASH树的根节点HASH值,然后用户可以使用从服务器获得的公钥验证返回的数字签名是否跟计算出的根节点HASH值对应。因为单向HASH方法自身的特性,如果攻击者对于返回的结果或者是验证对象做任何的修改都会导致客户端的验证出错。

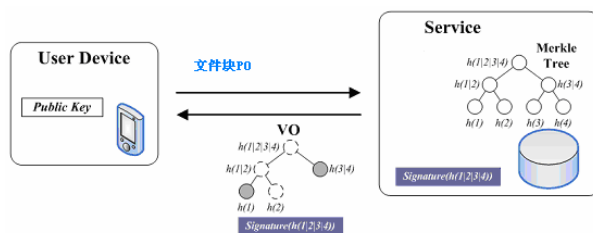


图3 基于Merkle HASH树文件校验过程

(1) User Device(简称U)向Service(简称S)发送下载文件P的请求;

(2) S依据配置信息,处理U的请求,将P文件的数字签名VO和文件描述信息Info,如文件大小,文件分块大小等发送回U;

(3) U收到S响应的信息,生成下载描述文件,保存下载文件P的VO和Info;

(4) U多线程发送下载文件块p0, p1, p2, p3请求;

(5) S多线程响应U的请求,多线程响应p0文件块数据和HASH(p0)值;p1文件块数据和HASH(p1)值;p2文件块数据和HASH(p2)值;p3文件块数据和HASH(p3)值;

(6) U接收各个文件块数据和HASH值,计算得到根节点的HASH;

(7) U 计算根节点的数字签名  $VO'$ ，比较 (3) 中的  $VO = VO'$ ？

(8) 如果相等，下载到的文件 P 是完整的，否则 P 是不完整的。

由此可得，Merkle Hash 树的方法之只需要计算一系列的 HASH 值而不是数字签名，而单向 HASH 函数的计算机代价比数字签名的代价小很多，所以 Merkle Hash 树的方法对简单数字签名的方法在效率上有很大改进。

采用 Merkle 哈希树校验方式能够极大地减小校验文件的尺寸，从而有利于缓解服务器集中下载瓶颈，Merkle 哈希树校验方式与分布式哈希表技术势必能够帮助无线应用下载协议进一步克服自身的非结构化缺陷，取得更大的应用发展。

### 3 基于Merkle树移动平台下载设计和实现

如图 4 所示下载过程时序图。对于手机下载业务而言，必须考虑以下几个方面的问题：如何找到下载内容；如何确保下载的内容是正确的、完整的；如何避免下载的内容未经许可的转发；如何保证正确的计费。本文研究下载内容的完整性问题，即如何保证到达终端的已下载内容是完整的。下载内容的完整性校验正是基于该问题的解决方案。根据完整性校验结果，用户可准确的获知自己得到的内容是否是完整的，安全的。从而，确保为用户提供安全，可信的内容下载服务。

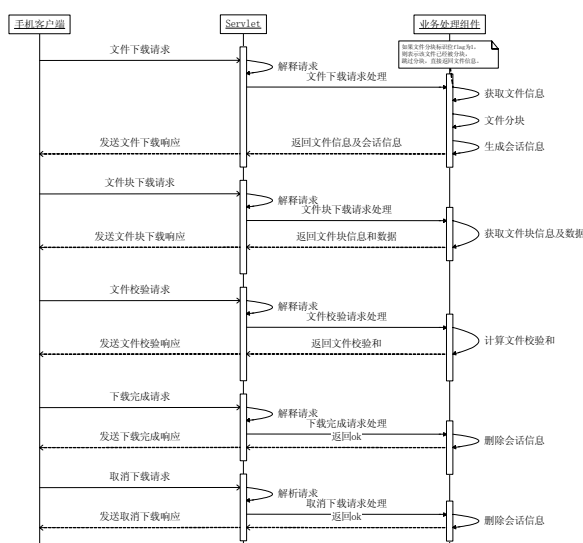


图 4 移动平台下载业务时序图

手机终端发送下载文件请求，服务器端接收到请求，获取资源信息，生成下载描述文件响应手机终端，终端用户多线程，分块下载资源，分块的信息、断点下载信息全部记录在下载描述文件中。

传统的多线程、分块下载过程中，文件的完整性校验是对每个分块文件都进行校验值计算，并在每块文件的传输过程中，额外携带文件校验值，此外在客户端需要保存每个分块文件的校验值。如需对文件来源证实的话，还要计算每块文件的数字签名。若手机终端作为下载客户端，考虑到客户端自身特征，传统方法俨然制约下载性能。例如：下载 100k 的文件，每块文件 4k，分块数为 25 个。传统的校验方法，服务器端每次传送到客户端最终文件大小为  $(4k+128bit)$ ，除最后一次传送整个文件的 MD5 值。假设每块文件的下载都是完整的，那么文件下载完整，需在网络上传递的数据量为  $(4k+128bit)*25+128bit$ ；此外，还要考虑客户端每次计算 MD5 值花费的时间  $T(MD5, \text{客户端计算性能 } p)$  和校验值在客户端的存储占用空间  $M(128bit*25+128bit)$ 。与传统文件校验方式相比，基于 Merkle 树的文件完整性校验，从服务器端向客户端网络传送数据量到客户端的校验值计算量，校验值存储所占空间都约为传统方法的 1/2。采用 Merkle 哈希树校验方式能够提高校验性能，这是因为，一旦采用这种方式来校验数据块，那么便没有必要在校验文件中保存所有数据块的校验值了，其中仅需保存一个 20 字节的 MD5 哈希值即可，即整个下载文件的根哈希值。

### 4 结束语

在本文阐述了基于 Merkle 树实现移动应用平台下载业务中，文件完整性校验策略，使用基于数字签名的方法可以保证服务器返回数据的正确性和完整性。但是考虑到实际的业务需求，该策略可以改进。例如：下载内容很大，考虑到网络传输的限制，下载文件分块数很多，导致原有的 Merkle 树的高度值很大，直接影响到客户端进行 HASH 值计算的次数增多。如果一个下载过程，占用大量计算机，应用其他应用的运行，这是不能允许的。改进的方法考虑不用二叉树，而使用效率更高的 B/B+ 树作为 Merkle 树的实现数据结构。

(下转第 162 页)

2008 年 2 月至 2008 年 7 月的交通事故数据挖掘, 经过数据筛选和清洗后, 得到 1054 条有效记录, 将其作为挖掘数据集, 选取其中最显著刻画交通事故发生规律的 4 个字段: (月份, 时刻, 发生地横坐标, 发生地纵坐标) 进行聚类挖掘。算法参数设置:  $n=1054$ , 数据维数:  $d=4$ , 迭代次数:  $N=60$ , 迭代阈值:  $\text{Delt}=0.001$ , 聚类有效性函数选取常用的 Xie-Beni 函数。对于聚类数  $C$ , 可先给定  $C$  的一个大致范围为 2~10, 然后根据最小的聚类有效性值确定最佳聚类数  $C$  的值。

算法最终挖掘结果为: 最佳聚类数  $C=5$ , 聚类中心为:  $C_1(2.358, 18.235, 232.247, 733.531)$ ,  $C_2(2.816, 17.943, 657.792, 1031.506)$ ,  $C_3(5.134, 12.134, 658.790, 267.714)$ ,  $C_4(7.463, 12.025, 1020.656, 1930.978)$ ,  $C_5(7.562, 11.841, 431.326, 1642.579)$ 。该结果表明, 2008 年 2 月至 7 月间, 该市交通事故多发生于 2 月, 5 月和 7 月的中午和傍晚时段, 并集中发生于  $(232.247, 733.531)$ ,  $(657.792, 1031.506)$ ,  $(658.790, 267.714)$ ,  $(1020.656, 1930.978)$ ,  $(431.326, 1642.579)$  五个地理位置附近处。

据此分析, 如果在上述时段和地理位置及其附近处有针对性的加强交通警力部署和交通标牌设施的设

置, 可以有效地减少交通事故的发生, 挖掘结果也表明, 改进后的聚类挖掘算法是很有效的。

## 5 结束语

数据挖掘技术应用于公安业务数据, 对于指导公安实战工作具有一定的指导意义, 并具有较好的应用前景, 可以利用其构建具有智能性和主动性的信息平台, 以提升信息化办案水平, 本文的后续工作将会在这方面作一些相关的研究。另外, 挖掘前业务数据的预处理及算法的较优参数选取将会在很大程度上影响挖掘的质量, 后续工作也会在这方面作进一步的探究。

## 参考文献

- 1 薛京生, 孙济洲, 孙宇, 何宏. 基于应急事件响应的模糊聚类分析算法. 计算机工程, 2006, 32(1): 201-202.
- 2 高新波. 模糊聚类分析及其应用. 西安: 西安电子科技大学出版社, 2004. 49-55.
- 3 Zadeh LA. Fuzzy sets. Information and Control. 1965, (8): 338-353.
- 4 杨兴春, 李进. 一种基于多种群隔代融合的遗传算法. 计算机与数字工程, 2008, 36(5): 30-32.

(上接第 176 页)

## 参考文献

- 1 Merkle R C. A certified digital signature. Advances in Cryptology-CRYPTO'89. Berlin: Springer-Verlag, 1989, 218-238.
- 2 Habib A, Xu D, Mikhail A, et al. A Treebased Forward Disgest Protocol to Verify Data Integrity in Distributed Media Streaming. Proc. of Eurocrypt'01, Lecture Notes in Computer Science, Vol.2007.
- 3 中国商用密码认证体系结构研究. 课题组. 数字证书应用技术指南. 北京: 电子工业出版社, 2008
- 4 王东滨, 方滨兴, 云晓春. 基于文件完整性校验的入侵检测及恢复技术的研究. 计算机工程与应用, 2003, 31: 156-157.
- 5 张雪旺, 唐贤纶. MD5 算法及其在文件系统完整性保护中的应用. 计算机应用, 2003, 12(23).
- 6 Merkle RC. A certified digital signature. Proc. on Advances in Cryptology, 218-238.