

大数据场景下的云存储技术与应用

Cloud Storage Technology and Applications for Big Data

陈杰/CHEN Jie

(中兴通讯股份有限公司, 广东 深圳 518057)
(ZTE Corporation, Shenzhen 518057, China)

摘要:文章认为随着大数据应用规模的扩大,新业务环境和场景对海量云存储需求的迫切,云存储平台需要打破原有的框架,改变组网和管理方式,以满足新的业务需求。文章分析了各种场景,提出了云存储的需求及关键技术等。文章指出大数据需求促进了云存储的发展,而云存储的发展则带动了新的业务应用。

关键词:大数据;云存储;安全

Abstract: With the expansion of big data applications, mass cloud storage has become a more important requirement. To meet service demands, cloud storage needs a new framework and new networking and management methods. In this paper, we discuss the various scenarios of cloud storage and discuss the demands and key technology of cloud storage. Big data requirements promote the development of cloud storage, and cloud storage development creates new service applications.

Key words: big data; cloud storage; safety

中图分类号:TP393 文献标志码:A 文章编号:1009-6868(2012)06-0047-05

也将迎来快速发展的机会^[1]。

1.2 大数据的类型

在经历了20世纪的计算浪潮和网络浪潮之后,信息存储技术已经发展成为信息领域的三大支撑技术之一。随着云计算、物联网等信息技术的飞速发展,异构数据源越来越多,数据信息量在飞速增长,数据的类型也复杂多样,不仅使得信息系统规模日益庞大,也导致海量非结构化数据管理复杂、异构数据存储利用率低下、资源不易扩展等问题。

海量非确定性异构数据产生的原因复杂多样,在应用中也具有新的特点:随着各种应用规模及领域的扩大,数据量会呈现爆炸性增长及海量数据存储的趋势和特点;在非确定数据的典型应用中数据源很多,数据种类也繁多,数据资源具有异构性特点;数据还呈现数据块大小、数据类型和数据访问方式等不确定的特点;云计算、物联网等应用的不断丰富,数据产生、应用、访问方式十分复杂,还使得数据具有时效性和空间性,高频度访问和高并发的特点^[2]。

1.3 大数据对存储的需求

非确定数据应用中的海量数据对数据的存储体系结构带来了很大

1 大数据应用场景与需求

1.1 大数据的发展

随着互联网、移动互联网、物联网的发展,大数据逐渐成为发展的趋势。数据不仅仅正变得更加可用,同时也变得更易被计算机所理解。大数据发展趋势中所增加的大部分数据都来源于物联网世界中商品、物流信息,企业内部经营交易信息,互联网世界中人与人交互信息、位置信息等。

2011年分析调研机构IDC发布的研究报告《从混沌中提取价值》显示:2011年全球被创建和被复制的数据总量为1.8 ZB,和2010年同期相比,这一数据上涨了超过1 ZB,而且这些数据大部分是非结构化的数据。预测到2020年,全球数据量暴增44倍(相比2009年),总量会达到35 ZB。根据图灵奖获得者Jim Gray

提出的数据增长的经验定律,网络环境下每18个月产生的数据量等于有史以来数据量之和,因此海量数据处理的需求会变得非常普遍。

2012年,《纽约时报》称大数据时代已经降临,决策行为将基于数据和分析,而并非基于经验和直觉。这不是简单的数据增多的问题,而是全新的问题。旨在从互联网时代非结构化数据的庞大宝藏中获取知识和洞察力的计算机工具也正在迅速发展。大数据技术领域的竞争,事关国家安全和未来,国家层面的竞争力将部分体现为一国拥有数据的规模、活性以及解释、运用的能力。2012年3月29日,奥巴马政府投资2亿美元启动大数据研究与开发计划,是大数据技术从商业行为上升到国家意志的分水岭。预计欧盟、中国等大型经济体很快也会出台相应倾斜性政策,大数据相关产业链公司

的挑战。首先,海量数据的组织必然采用分布式数据组织与管理策略,这需要实现适合于非确定数据应用的(元)数据和数据组织方式;其次,由于海量数据是通过持续增长积累而成,而积累的过程需要很长的时间,因此需要存储支持可保证规模与性能同时扩展的存储组织模式以及相应的索引机制。

针对海量不确定性数据,使用基于传统的信息存储结构和对象查询方法的实际运行效率呈现下降趋势,因此必须采用新的元数据组织结构和查询方法来提高效率,为用户提供高性能的多并发数据查询服务。

由于在分布式环境中,数据源分布在不同的网络结点,这就存在网络传输性能低的问题。而各个数据源有很强的自治性,它们可以自治地改变自身的结构和更新数据,这就会给数据集成系统的一致性带来了困难。由于数据存在非确定性,针对海量非确定性异构数据的集成工作将变得更为复杂,可以采用分布式并行处理技术实现计算资源和存储资源的全局最优化的管理。

在信息化时代和全球经济竞争的新环境下,要做出一项决策,往往需要查询多个业务系统和外部系统,并进行大量的数据分析,工作量大且易出现人为差错,影响决策的质量。这就需要对海量非确定性异构进行整合集成,整合、集成后的数据必须保证一定的集成性、完整性、一致性和访问安全性。

数据的海量性、非确定性以及异构性为传统的数据挖掘算法提出了挑战。由于数据的异构、海量、分布性和决策控制的实时性,需要调整数据挖掘引擎的布局及多引擎的调度策略。结构化或者非结构化数据都涉及数据的存储、管理(索引、并发、一致性、查询等)等,这是因为用户对大数据使用方面的要求(对海量非结构化数据查询仍然要准确和快速),导致对数据逻辑结构和物理存储方

式的新要求。

随着大数据各种应用规模的扩大,数据量会呈现爆炸性增长及海量存储的趋势和特点。传统的数据存储系统达到了瓶颈,无法及时地完成各项运作任务。大数据的新特点对存储提出了新的挑战,为了适应大数据的发展,存储需要支持纵向无限扩容(存储扩容)和横向无限扩容(能力扩容),并以对象作为基本的存储形式,以提高系统的扩展性,降低系统维护复杂度。

2 云存储的关键技术

针对数据的飞速发展和数据安全要求的不断提高,如何建立安全、性价比高的存储成为业界的普遍需求。云存储成为首要选择,因为它能够根据所需容量大小对用户进行定制,用户不需要进行硬件的管理维护,缩减了用户成本和人力投入。而且云存储具有易扩容、易管理、价格低、数据安全、服务不中断等优点。

云存储是在云计算概念上延伸和发展出来的新概念,通过集群应用、网络技术活分布式文件系统等功能,将网络中大量各种不同类型的存储设备通过应用软件集合起来协调工作,共同对外提供数据存储和业务访问功能的一个系统。云存储是一个以数据存储和管理为核心的云计算系统。

大数据时代的来临对云存储提出了很多关键需求:

(1) 大规模级别存储系统的构建

随着数据的爆炸性增长,存储的规模越来越大。2012年云存储的建设规模是几十PB级,存储的文件数或者对象数是几十亿。到2013或2014年,就会有百PB级和EB级的需求,过几年将会增长到ZB级,文件数或对对象数也会超过百亿、千亿。

传统存储通常是在一个设备、一个机架或一个数据中心内完成资源组织管理,而当存储容量上升到EB级或ZB级后存储则很难在一个数据

中心内完成。大规模的存储需要跨数据中心,跨城市、省、甚至国家进行存储设备、存储数据、存储服务的组织和管理,并支持跨域的访问、备份、容灾等功能。同时大规模的存储要求存储提供不同等级的管理和服务权限,并按照区域、级别分配不同的权限。系统对资源的访问必须经过严格的权限控制。只有用户确认共享的资源才能被其他用户或业务进行访问,即使是被授权的访问也会根据不同的权限控制方式受到访问权限控制。

云存储就是将不同种类的存储设备协调起来进行工作。这些存储设备使用的存储介质也是多种多样的,而且随着技术的发展,设备种类和存储介质种类会越来越多,如何调度这些设备和存储介质协调工作,需要在云存储管理软件上考虑和优化,以保证组织好的资源被高效利用。

(2) 存储设备在线扩展和收缩

在存储设备的使用过程中,会遇到调整存储资源池的需求,这则要求存储资源池根据业务的需求增加或者减少存储设备。在调整的过程中,业务不能被中断,也不能使上层业务感受存储资源池的变化,同时被裁剪设备的数据要在较短的时间内在其他设备上恢复、备份,并在较短的时间内完成增加存储设备和原有存储你设备的数据均衡。

云存储系统要优化和调整数据组织和管理方法,即使存储规模增加后,性能要随之线性增加。数据变得庞大后,元数据管理要考虑中心化或多节点方式,以降低元数据管理对整个系统读写性能的影响。对于热点数据支持自动的多副本复制,则会在多个存储节点提供读能力,以降低硬盘、网口、处理器对性能能力的限制。采用多级缓存技术,热点数据则会先读入智能加速卡,并由智能加速卡对外提供读服务,在写数据时,也是先写入到智能加速卡,由加速卡组织分发到存储设备上。

(3) 实现面向应用的专业化管理策略

实现面向应用的专业化管理策略呈现出一些特点:传统存储系统存储资源与应用独立,存储资源利用效率低;海量存储系统把资源进行了整合,但是针对所有应用都采用统一存储策略;在云存储系统中如何做到资源整合并且针对应用进行专业化的策略管理,根据应用的变化进行弹性配额管理是一个较大的挑战。

云存储必须提供基于容器的多层次租户/应用隔离技术:系统必须提供数据隔离功能,保证数据不被非法访问,并保证用户数据的隐私。云存储可以通过物理隔离与权限控制相结合,实现对数据的隔离。

- 提供以用户为单位的数据隔离:业务系统为每个用户创建独立的存储空间,业务系统根据用户标识和对应权限对用户存储空间的数据进行访问控制,这样可以避免未授权用户访问到其他用户的数据以及用户信息。

- 提供以业务为单位的数据隔离:在进行数据的存储和读取时,每个业务都必须拥有自己独立访问的权限,系统根据不同的业务将数据隔离,避免数据被未授权的业务访问。

- 提供以存储容器为单位的数据隔离:可以设定数据存储在指定的存储容器中,不同的存储容器有不同的访问授权。访问授权可以是基于用户的,也可以是基于业务的。

云存储提供基于系统或者应用的多种服务管理策略:提供压缩策略,用户可以根据文件类型与活跃度设置压缩条件;提供系统级、业务级和用户级的流量控制策略设置;提供系统级、业务级和用户级的数据分片设置,业务可以设置对象存取的分片大小、分片存储区域(跨盘、跨节点、跨区域),同一对象的各分片可并发存取;提供系统级、业务级和用户级的热点对象设置,业务可以依据对象的访问活跃度设置热点对象和热点

对象的存取方式;提供系统级、业务级和用户级的文件归档设置,业务可以设定归档区域和归档条件,包括对象活跃度、容器活跃度、文件类型等;提供系统级、业务级和用户级的隐私性保护,对象及其元数据必须归属用户,非用户授权任何用户不允许访问;提供系统级、业务级和用户级的重复数据删除策略设置,可以按命名空间、存储容器、存储区域、文件类型、执行时间等设定重复数据删除的策略,并可以查询重复数据删除的操作记录和效果分析。

云存储要提供弹性资源伸缩和共享:系统要支持根据业务使用情况自动的增加和缩减存储空间,同时利用重复数据删除技术,提升存储资源的利用率。

(4) 系统全局自动负载均衡

在云存储的系统中,物理存储主机节点规模从几万到几十万,多为数据密集型应用,比如每天亿次级别的网络搜索访问。面对超大规模的数据请求和节点数量,如何高效进行节点负载均衡,如何发挥空闲节点的作用是保障高水平服务质量,提高系统运行效能一个较大的挑战。

云存储系统要求云存储具备基于服务质量(QoS)的多层次自动负载均衡与调度功能。

- 实现基于请求类别及前端节点负载进行的均衡和差异化调度:系统参照业务诉求和区域的QoS信息(存储总容量、总的IO吞吐能力、当前系统繁忙程度等)为业务选择最合适的区域归属,如对于IO要求较高的应用可以放到IO吞吐能力较强的区域里。

- 实现基于请求类别及数据分布进行的均衡和差异化调度:系统必须提供对多组云存储系统之间的动态调度能力,并根据区域内每组系统的IO繁忙程度,将业务的访问请求尽可能发送给那些IO负荷不重的组,以实现组间IO的负载均衡。

- 实现数据中心之间的负载均衡

度,均衡各中心利用率:通过调度服务和资源策略实现资源的跨域整合,还可以通过访问重定向和数据迁移等多种技术手段对外提供统一的存储资源服务,并对用户屏蔽资源的具体位置信息并自动实现就近访问。

此外,云存储设备还采用数据压缩技术构建分布式缓存系统,提高给定缓存加载的数据量,提升系统性能。同时,在广域网数据传输前进行重复数据检测,相同数据只传一份,就可以实现基于删冗的广域网数据传输加速。

(5) 云存储数据安全和数据保护

云存储系统需要支持的用户数量巨大,且存储了用户生活、工作、学习等各种类型的数据,具有私密性,另外对于数据的可靠性和完整要求也非常高。因此如何解决用户数据的共享和隐私保护之间的矛盾、用户数据的可靠保护和存储高效之间的矛盾是一个很大的挑战。

系统对资源的访问必须经过严格的权限控制,只有用户确认共享的资源才能被其他用户或业务进行访问,即使是被授权的访问也会根据不同的权限控制方式受到访问权限的相关控制。

云存储还需要具备以下的一些基本功能:

系统必须提供数据隔离功能,以保证数据不被非法访问并保证用户数据的隐私。通过物理隔离以及权限控制相结合,可以实现对数据的隔离。

系统必须提供信息加密的功能,防止用户信息被窃取。用户的关键信息,如登录密码和系统访问等其他鉴权信息,无论是传输时还是在存储时必须加密。

系统必须提供数据传输加密功能。数据的传输加密可以通过客户端软件的传输设置实现。用户设置采取加密通道传输,系统应当在重新登录后进行数据传输时使用Https通道进行数据的传输。

系统必须提供有效的硬盘保护形式,保证即使硬盘被窃取,非法用户也无法从硬盘中获取有效的用户数据。

系统必须支持数据加密存储,用户在使用客户端软件时可以选择对存储数据加密。采用数据加密存储的客户端软件在上传数据时对数据进行自动加密,在线备份获取加密数据后能够在客户端自动解密,而在云存储获取加密数据后必须手动解密。

系统必须将数据切片存储在不同的云存储节点和硬盘上,数据无法通过单个硬盘恢复。故障硬盘无需进行数据清除即可直接废弃,用户数据不会通过硬盘泄露。

(6) 云存储节能降耗

云存储系统规模巨大,且需要提供高质量的对外服务,传统的构建方法提高性能和能耗增长近似线性关系,需要用新系统的架构打破这种关系,解决系统性能和能耗的矛盾。

- 建立不同级别的存储池:高写高读、高写低读、低写高读、低写低读等,并按照业务模型进行资源分配。对于不同的存储资源池采用不同的策略:如重删、压缩、按需分配等,提高存储的利用率,从而降低系统能耗以及二氧化碳的排放量。

- 监控业务访问量,控制系统内部部分设备可以轮流进行休眠及CPU降频。同时考虑使用低功耗的处理器。

- 控制业务流向能源供应充分或者消耗能源少的存储资源池。

3 云存储的应用

从2010年开始,云存储的运用越来越广泛,在互联网、平安城市、视频制作、数字传媒、家庭娱乐、个人网盘等方面都有很多应用。

3.1 视频监控应用

随着城市的现代化建设和经济的快速发展,构建和谐社会的必要性与日俱增,每个城市都在努力打造

平安城市。平安城市等大规模高清视频监控系统中的主要问题就是如何处理庞大的高清视频数据,因此就必须从视频的采集、编解码、传输、实时监控、录像回放等环节全面支持大规模高清。这样就给高清监控系统带来了一系列问题:网络带宽紧张、存储空间庞大、对性能的要求成倍增长、系统扩容升级压力等。如果采用1080P的高清视频监控,即使使用具有高压缩比的H.264编解码技术压缩高清视频,输出码率也将可达到6 Mbit/s,那么每台摄像机每天将大约产生50 GB左右的数据量,一个月就是1.5 TB左右的存储量。一个城市有上千个摄像头,一个月的数据量就达到PB级以上。面对PB级的海量网络存储需求,传统的开放系统的直连式存储(DAS)和网络存储技术(NAS)在容量和性能的扩展上存在瓶颈,已经不能满足对高清视频的存储。

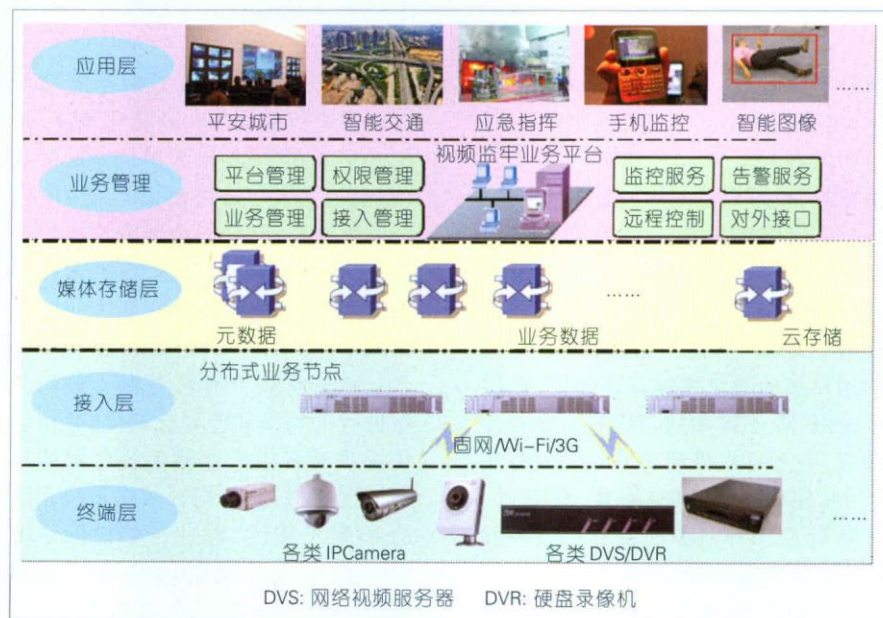
远程云端用户可以通过以太网访问云存储中的视频,如图1所示,云存储提供多种应用接口给视频监控系统中的应用、管理和使用。采用基于浏览器的方法,云端用户不需要安装任何播放以及管理软件,就可以远程对监控录像进行回放、视频的

再分析等。

云存储为实现大规模高清视频的存储和处理提供了一种新的解决方案。云存储突破了传统存储方式的性能和容量瓶颈,实现了性能和容量的线性扩展,让海量数据存储成为可能。同时云存储可以实现存储完全虚拟化,所有设备对云端用户完全透明,任何云端、任何被授权用户都可以通过一根网线与云存储连接,从而让用户拥有相当于整片云的存储能力。由于各个监控区域地理范围分布广阔,监控点的数据巨大,采用云存储系统,便于分布式管理和随时扩容。云存储由元数据管理节点和大量存储节点组成。元数据节点负责存储虚拟化、资源管理以及存储数据的命名空间、存取控制、存放位置等信息,是云存储的核心部分。云端不直接通过元数据节点读取数据,而是从元数据节点获取视频存储的位置信息后,直接和存储节点进行读写操作。将控制部分和业务数据部分分离,有助于提高系统的可扩展性和数据处理的读写带宽。

3.2 互联网应用

一个用户拥有多个终端的现象



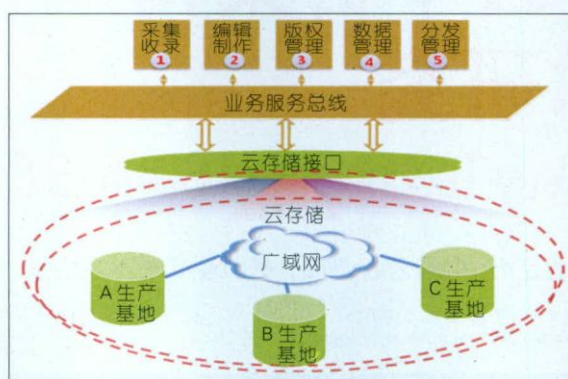
▲图1 视频监控中云存储的应用

越来越普遍,用户的同步和分享需求则会变得越来越强烈,同时移动终端以及移动互联网的快速发展均为个人云存储的发展提供了良好的发展环境。

互联网公司的纷纷加入云存储的研发和应用队伍中,并将互联网产品中需要存储的个人信息与云存储应用绑定,加快了用户的接收速度。基于互联网提供的云存储平台,各类消费电子产品实现了前所未有的互联互通,包括文字、图片、音乐、视频、应用等在内的数字内容开始实现跨越时间和空间的自由流通,为社交网络中的用户提供了良好的交流元素。这些数字内容目前越来越向诸如 Dropbox、Box.net、115 网盘等一些专业云存储服务商集中。

要求云存储支持 10 PB 以上容量空间,提供网络文件系统(NFS)等文件或对象访问接口。如图 2 所示,云

管理的复杂性,增强了系统的灵活性和可扩展性,满足了海量视频数据的存储需求。图 3 所示说明了视频数



▲图3 云存储跨地域组网

据可以跨地域存储和调用。

云存储系统提供高速的读写接口,满足采编等应用的视频加工需求,实现对原始素材、成品节目、再加工节目等不同类型节目的分层存储和分类管理。

现第四屏、机顶盒等的互操作。

4 结束语

随着互联网、移动互联网、物联网的发展,大数据逐渐成为发展的趋势,数据产生的原因复杂多样,在应用中也具有新的特点。随着各种应用规模的扩大,数据量会呈现爆炸性增长的趋势及海量数据存储的特点。新业务环境和应用场景对海量云存储需求越来越迫切,这需要海量存储平台打破原

有的框架,改变组网和管理方式,满足业务需求。文章分析了各种场景,提出了对云存储的需求、关键技术和指标。文章认为云存储的应用越来越广泛,云存储技术的发展促进了业务融合,并衍生出新的业务应用。总之,大数据的需求促进了云存储的发展,云存储的发展带动了新的一系列业务应用。

5 参考文献

- [1] 王伟,柯尊友.云存储的进化:云存储解决方案[J].中兴通讯技术(简讯):2012(8):18-19.
- [2] 李群.国内个人云存储应用风生水起发展迅猛[EB/OL].<http://net.chinabyte.com/115/12393115.shtml>

收稿日期:2012-10-12



▲图2 互联网中云存储提供的接口方式

存储提供满足多种应用需求的接口。由于互联网对云存储的访问速度要求相对较低,关注建设成本,因此需要云存储提供纠删码方式的数据保护、提供压缩、重删等多种策略。并可以基于容器或者用户为单位对外提供租赁服务,按照空间大小、访问流量进行收费。

3.3 视频编排应用

云存储系统通过对物理存储设备的虚拟化管理,实现视频数据与物理存储位置无关性,降低了视频数据

3.4 家庭娱乐应用

实现家庭网络内各终端通过云存储访问代理与云存储平台交互,则可以完成媒体文件的上传和在线播放,并实现家庭网络内各终端信息在云存储上共享和备份。

基于家庭的云存储可以实现家庭内多终端间的多媒体资源共享及多屏互动、多终端交互控制,优化

用户感知。

图 4 所示为家庭媒体云。通过家庭媒体共享和交互控制功能,可以实



▲图4 家庭媒体云

作者简介



陈杰,南京邮电大学通信专业、纽约大学计算机专业双硕士学位毕业;曾任中兴半导体有限公司开发部主任、AT&T公司贝尔实验室高级研究员与研究部主任、中兴通讯美国分公司总裁、中兴通讯网络事业部总经理,自2003年以来,任中兴软创董事长、中兴通讯高级副总裁;1992年获得国家科学技术进步奖一等奖,2007、2008年分别获得国家科学技术进步奖二等奖,2012年获得深圳市高层次专业领军人才称号。