为一款成熟的企业级分布式数据库,基于普通PC服务器,就能够做到传统高端硬件环境下的数据可靠性和服务可用性,而且还能做得更好!跟我们一起看看OceanBase的技术秘诀吧!

「作者介绍

「中者介绍

「中本介绍

「中本の表別

「中本の

责OceanBase数据库的外部

客户推广及技术支持工作。

前言

在「不可靠」硬件上,分布式数据库

如何保证数据可靠性和服务可用性?

OB君: "数据不能丢,服务不能停",可以说这句

话道出了用户对数据库的核心能力的要求。然

而, 传统的商业数据库必须依赖高可靠的硬件才

能实现数据可靠性和服务可用性。OceanBase作

这句话也道出了用户对数据库的一个核心能力要求:除了功能完善、使用方便之外,还要绝对安全、足够健壮。可以说,为了满足这两个看似简单的要求,在数据库领域诞生了大量的技术和论文,也让无数人绞尽了脑汁。

在传统的商业数据库产品(如Oracle、DB2)中,虽然也有一些行之有效的软件技术(如Redo Log、主从热备技术等)用来提高数据可靠性和服务可用性,但整体来说对硬件的稳定性有很强的依赖。而传统的企业级服务器(如IBM 的Mainframe、AS400、Power等)和EMC、IBM等厂商的高端存储产品,能够很好的保证硬件的稳定性,因此也就成为了Oracle为代表的传统数据库产品的理想平台,这也就是"IOE"一词的由来。可以说,I和E的重要职责就是保障O的稳定运行。

乔国治(鸷腾)

Part 1

说到数据可靠性和服务可用性,在数据库领域真是老生

常谈的话题,可以说从数据库诞生之日起就如影随形。

如果要用一句话来概括数据库对数据可靠性和服务可用

性的要求,可以借用OceanBase数据库创始人阳振坤

老师的一句话:"数据不能丢,服务不能停"。可以说,

EMC存储
Oracle 数据库
AIX、RS6000、OS400
IBM 服务器:
小型机(i和p系列)、大型机(z系列)

不知不觉间,IT世界进入PC服务器的时代。和传统架

构相比,PC服务器能降低成本并带来扩展性上的便

利、逐渐成为以互联网为代表的众多企业用户的首选。

但PC作为服务器也为用户带来了一个棘手的问题、那

就是硬件(服务器、内置磁盘、网卡等)的可靠性明显

下降了。虽然也有一些其它机制(比如RAID、外挂磁

盘阵列等)可以用来改善这种情况,但不能从根本上解

决问题。

2)

4)

主从热备

存储层数据校验

3) 备份/恢复

近些年来新兴的数据库产品,尤其是分布式数据库,几 乎无一例外地采用了PC服务器架构,**那在这种相对不** 稳定、不可靠的硬件条件下,分布式数据库如何保证数 据可靠性和服务可用性呢? 对习惯了传统高端硬件稳定 性的企业用户来说,这个问题的答案将直接决定业务要 承受多大的风险。 针对这个问题, OceanBase给出了非常明确的答复: 在PC服务器架构下,OceanBase不但能够满足传统高 端硬件环境下的数据可靠性和服务可用性,而且还能做 得更好!可以说,新技术趋势(PC架构)所带来的问 题(硬件可靠性降低),反而倒逼了技术自身的演进, 大家开始认真思考如何在"不可靠"硬件环境下保证数据 的"可靠性",并进一步保证服务的"可用性"。最终的结 果,是更多的从软件层面引入保障机制,来弥补硬件环 境的不足。后文将从多个方面来阐述OceanBase具体 是如何做的。 Part 2 OceanBase如何保证数据可靠性 在传统数据库中,有几种常用的手段来保证数据可靠 性: 1) Redo Log

都无法做到完美(即RPO=0),也就是无法保证数据完全无损,下面我们简单分析一下原因。

首先看Redo Log。采用Write-Ahead-Log(WAL)模式的Redo Log可以保证数据库中已提交的数据不会丢失,如果已提交的数据还在内存中就发生了宕机等意外,利用Redo Log可以恢复这些还未持久化的数据。但这里有一个前提,就是Redo Log自身必须绝对可靠,如果Redo Log所在的存储发生损坏,那么这一前提便不复存在。

因此,Redo Log所带来的数据可靠性其实取决于硬件的可靠性,说到底还是要依赖高端硬件。此外,如果是已经持久化的数据遇到了硬件损坏(比如"坏页"问题),并且对应的Redo Log已经被覆盖,那么Redo Log也无能为力了。

有了主从热备技术(比如Oracle DataGuard,IBM

这些技术可以从很大程度上提高数据的可靠性,但似乎

DB2 HADR等)之后,Redo Log可以写两份甚至多份了,即使主节点遇到硬件故障,仍然可以用备节点的Redo Log来恢复数据,看上去可以做到RPO=0了。但其实不然,虽然主从热备技术通常都提供"数据强同步"的手段(比如Oracle Data Guard中的Max-Protection模式,IBM DB2 HADR中的Sync模式)来保证RPO=0,但实际系统中几乎没有人采用。为什么呢?原因很简单,这种模式下一旦备节点发生故障,或者主备之间的网络发生故障,那么主节点的正常交易就会受拖累;换句话说,备节点不但没有提高整体稳定性,反而降低了整体稳定性,得不偿失。所以,在实际生产系统中部署过主从热备的朋友都知道,几乎没有人会采用

数据强一致模式,也就无法做到RPO=0。

件的可靠性。

除了Redo Log以外,我们还有用来保底的"备份/恢

复"手段。但备份/恢复机制在时效性上有明显缺陷、只

能最为应对最差情况的"最后一发子弹",不到万不得已

绝不使用。而且,备份文件的可靠性最终依然是依赖硬

最后,数据库的存储层通常都会有数据校验机制,用来

检测存储层的"静默错误"和潜在的软件错误。但数据校

验机制更多的是用来在事后发现错误,无法预防或者解

决错误,后面我们会介绍OceanBase如何将这种技术

综合上面的分析,可以看到传统数据库最终还是要依赖

和其它机制结合起来,以提高整体的数据可靠性。

高端硬件的可靠性来保证数据的可靠性,只依靠数据库 自身的能力是无法保证RPO=0的。 那么在PC服务器的环境里,硬件的可靠性明显不足, OceanBase这样的分布式数据库又是如何来保证数据 可靠性呢? 前面提到过,我们更多的是在软件层面引入保障机制, 比如以Paxos(以及后面衍生出来的Raft)为代表 的"分布式一致性协议"。OceanBase充分利用了Paxos 协 议 , 并 将 Paxos 协 议 和 传 统 的 WAL 机 制 结 合 起 来,每一次Redo Log落盘时,都会以强一致方式同步 到Paxos组中多数派(leader+若干follower)副本的 磁盘中,这样做有两个好处: 1)在Paxos组中任意少数派副本发生故障的情况下, 剩下的多数派副本都能保证有最新的Redo Log,因此 就 能 避 免 个 别 硬 件 故 障 带 来 的 数 据 损 失 , 保 证 RPO=0. 2) Paxos协议中的数据强一致是针对"多数派"副本而

言,而不像主从热备那样要求"所有"副本的数据都保持

强一致。如果Paxos组中有少数派follower副本发生故

障,剩下的多数派副本(leader+若干follower)之间

的数据强一致完全不受影响,这就解决了前面提到的问

题: 主从热备模式下备副本故障拖累主副本的可用性。

综合以上两点,OceanBase利用Paxos协议可以保证

RPO=0,且不必担心应用的性能会受到影响,这也是

OceanBase和传统数据库在数据可靠性方面最显著的

前面提到过的"数据校验"机制,在众多工程实践中已经

被证实为一种行之有效的机制。但是,如果要在采用

Paxos协议的分布式数据库中实施数据校验机制,情况

不同点。

将更加复杂:除了要关注单个物理节点的"静默错误", 还要保证多个副本之间数据的一致性。 假设网络传输出了问题,导致错误的数据从leader副本 发送到了follower副本;如果follower副本在不知情的 情况下,将这些错误的数据当成正确的数据存储起来, 一旦leader副本发生问题,而有错误数据的follower副 本又接管了服务,那么久直接造成了用户的数据损失。 OceanBase也在存储层引入了数据校验机制,但我们 加入了更多的技术手段以避免上述问题,大致包含以下 内容: 1) Redo Log的数据校验 首先, Redo Log在落盘的时候会加上数据校验信息, 用来应对可能发生的磁盘静默错误。此外,为了保证一 个Paxos组中多个副本之间Redo Log的一致性, Redo Log在leader发送和follower接受时都会检查数据校验 信息,避免网络传输问题导致的数据错误。 2)数据盘上的校验信息

和Redo Log类似,数据盘上的数据也会有校验信息以

应对磁盘静默错误。但由于OceanBase是通过Redo

Log实现Paxos组中多个副本之间的数据同步,数据盘

上的数据并不会通过网络传输在多个副本间同步,因此

刚刚提到了,OceanBase不会对数据盘上的数据做副

本间的"实时"校验,但我们还是会在一些特定的检查

点,对多个副本之间的数据盘做一致性检查。这个检查

点选在了OceanBase的"每日合并"点,主要的原因,

是每日合并动作本身就要对大量数据做归并和重新写

入,刚好可以利用这个时机做数据的一致性检查。通过

这个检查,进一步在存储层确保了多个副本之间的数据

对于有关联关系的数据对象,OceanBase会做额外的

检查以保证它们之间的数据一致性。比较典型的例子就

是索引和它的数据表,OceanBase会在一些特定的检

查点(如每日合并点)做索引和数据表之间的一致性检

不需要副本间的实时校验。

3) 副本间的检查点一致性校验

一致性,提高了数据可靠性。

查,进一步提高数据可靠性。

5) 定期做数据校验信息检查

4)数据表和索引表之间的数据一致性校验

上面提到的一些数据校验措施(比如Redo Log和数据 盘上的数据校验信息),主要目的还是在数据中埋入校 验信息。但光有校验信息是不够的, 还要能够利用校验 信息及时发现磁盘的静默错误, 否则就只能等到访问数 据的时候才能发现错误,为时已晚。为了应对这个问 题,OceanBase在后台有定期的检查任务,在不影响 在线业务的前提下,利用数据校验信息主动检查磁盘静 默错误,一旦发现错误会及时通知用户,尽快采取补救 措施。 最后、OceanBase也和传统数据库一样提供完善的备 份/恢复机制,包括全量备份功能和增量备份功能。而 且OceanBase的增量备份是以不间断的后台daemon任 务形式持续进行,完全不影响在线业务,降低了运维操 作的复杂度。不过从分布式数据库的运行实践来看,在 实际系统中极少发生Paxos组中多数派副本同时毁坏的 情况,因此基本不会真正用到备份来恢复数据。 Part 3 OceanBase如何保证服务可用性

在数据可靠性有保证的前提下,服务可用性就成为了另

一个焦点:如果某个服务节点发生了故障,用户不但希

望数据不丢(RPO=0),而且希望服务能够尽快恢复

在历史上,在传统数据库的高可用能力不足时,有很多

种高可用方案是结合了硬件和操作系统的高可用能力,

不过这些方案通常在架构上过于复杂,而且无法在数据

库层面保证数据的一致性。随着数据库内置的高可用能

力逐渐完善,用户也转向了数据库自带的高可用方案,

典型代表就是"主从热备"技术, 比如 Oracle Data

但是,主从热备技术在高可用上有一个很难解决的问

题: 当主节点故障的时候, 如何让备用节点快速接管服

务。如果让备用节点判断主节点的状态,并且在主节点

故障时"自动"接管服务,那么就会面临一个致命的问

题,就是"脑裂(Split-brain)":备用节点和主节点同

时提供服务了,两个节点间的数据将再也无法保持一

为了避免脑裂问题,数据库的主从热备机制都不会提

供"备库自动接管服务"的能力,备库的接管动作要引入

人工决策的流程,需要人为在备用节点发起服务接管的

动作。这样一来,RTO就会达到数十分钟甚至以小时

计。为了避免脑裂问题,同时又减小RTO,有些数据库

产品引入了自动failover的机制,比如Oracle数据库的

Fast-Start Failover(FSFO);但由于引入了额外的组

件或者服务,整个解决方案的复杂度明显增大,而且这

些组件或者服务自身的高可用又成了新的问题。总体来

说,由于主从热备的本质没有发生改变,只是额外引入

了第三方仲裁者的角色,所以这种方案并没有解决根本

前文已经提到, OceanBase利用了Paxos协议中的多

数派共识机制来保证数据的可靠性, 在高可用方面,

OceanBase是利用了同样的机制。首先,根据Paxos协

议,在任一时刻只有多数派副本达成一致时,才能推选

一个leader,其余的少数派副本则不具备推选leader的

其次,如果正在提供服务的leader副本遇到故障而无法

继续提供服务,只要其余的follower副本满足多数派并

且达成一致,他们就可以推选一个新的leader来接管服

务,而正在提供服务的leader自己无法满足多数派条

件,将自动失去leader的资格。因此,我们可以看到

1) 从理论上保证了任一时刻至多有一个leader,彻底

2) 由于不再担心脑裂,当leader故障而无法提供服务

时,follower便可以自动触发选举来产生新的leader

这样一来,不但从根本上解决了脑裂的问题,还可以利

用自动重新选举大大缩短RTO,可以说完美解决了主从

Paxos协议在高可用方面有明显的优势:

并接管服务,全程无须人工介入。

杜绝了脑裂的情况。

问题,也不能100%保证避免脑裂。

(RTO越小越好)。

Guard、IBM DB2 HADR等。

致。

资格。

热备技术在高可用上所面临的难题。 当然,这里面还有一个很重要的因素,那就是leader出 现故障时,follower能在多长时间内感知到leader的故 障并推选出新的leader,这个时间直接决定了RTO的大 小。在OceanBase中,为了能够及时感知Paxos组中各 个副本(包括leader和follower)的状态,在各个副本 之间会有定期探活的消息。另一方面,探活机制虽然能 够检测到节点的故障,但是在网络不稳定的情况下,也 可能由于偶发的探活消息丢包而产生"误报(Falsealarm)"的情况。为了避免误报对系统稳定性带来的影 响,OceanBase也采取了很多应对措施: 1) 首先,探活消息的周期必须要合理。 如果周期太长就不能及时感知到节点故障,如果周期太 短就会增大误报的概率,而且也可能会影响性能。目前 在OceanBase中的探活周期为10秒左右,确保能及时 感知到节点故障,并且也不会频繁产生误报。 2) 其次,要能够容忍偶发性的消息丢包,减小误报的 概率。 具体来说,OceanBase不会由于一次探活的失败就认 定某个节点发生了故障,而是在连续多次尝试都失败后 才确认真正的节点故障。这就有效避免了偶发性消息丢 包所导致的误报。 3) 如果真的发生了误报,需要将影响范围降到最小。 OceanBase将Paxos组的粒度下沉到了表的分区一级,

也就是每一个分区都会有一个Paxos组,用来维护这个

分区的多个副本之间的leader-follower关系。如果由于

少量网络丢包导致"某个分区"的探活消息没有收到回

复,那么受影响的只是这个分区,同一台机器上的其它

分区会照常工作,这就有效地控制了问题的影响范围。

举个例子,如果某台机器出现间歇性故障(比如网卡或

者操作系统出了问题),导致这台机器频繁发生网络传

输故障,就会使这台机器上所有的leader副本持续受到

影响。这种情况下,OceanBase可以通过设置特定参

数限制这台机器暂时不参与leader选举,这样就有效地

起到了隔离作用,避免了局部故障对整个集群的可用性

在具备了上述这些处理机制后,OceanBase目前已经

OceanBase的容灾方案

4) 某些特殊情况的处理。

造成持续影响。

能做到最多10秒钟检测到服务节点异常,并在10~30秒 内完成服务的自动恢复。需要说明的是,具体的恢复时 间和遇到问题的机器个数、表的分区个数、故障类型 (机器硬件故障、网络设备故障等)都有密切的关系, 所以上面说的服务恢复时间只是作为一个参考值,在某

些特殊情况下也可能发生偏差。

Part 4

前面说到的内容,更多的还是从逻辑架构上介绍OceanBase如何实现数据可靠性和服务高可用。但实际应用中,除了逻辑架构之外,还必须考虑到系统部署时的物理分布情况,否则就无法充分利用Paxos协议所带来的优势。这就衍生出了机架容灾、同城机房容灾、异地城市容灾等诸多和容灾相关的概念。

之前在OceanBase微信公众号中已经介绍了常用的容灾方案,在这里就不再赘述了,有兴趣的朋友可以参考公众号里的文章:

• 一文详解OceanBase的高可用及容灾方案(下

Part 5 总结

到目前为止,我们已经介绍了OceanBase中关于保证数据可靠性和服务高可用的一些基本原理。原理看上去简单,但由于分布式数据库的网络复杂性以及PC硬件的不可靠,在实际使用中会发生形形色色的各种异常情况,任何一个细节处理不好都会造成严重后果。

经过在蚂蚁金服的大大小小、成百上千个系统中多年的打磨和锤炼后,OceanBase已经逐渐完善了众多的处理细节,形成了有效的机制来保证数据可靠性和服务可用性,而且,这套机制已经在众多实际系统中得到了很

理细节,形成了有效的机制来保证数据可靠性和服务可用性。而且,这套机制已经在众多实际系统中得到了很好的验证,尤其是支付宝和网商银行这种"线上数据就是全部,一个字节也不能丢"的金融级核心系统。

因此,今天OceanBase可以说已经突破了传统单点数据库在数据可靠性和服务高可用方面的限制,让用户的数据更安全,服务更稳定!

① 现场剪光缆!支付宝模拟自断一半服务器,26秒一切恢复正常。这一切的背后基于……
② 传统关系数据库高可用的缺失

● 大师专栏 | 如何在「不可靠」硬件上实现金融级

▼内容这么棒,还不赶紧扫码关注一下!▼

高可用?

OceanBase

金融级分布式关系数据库系统

长按关注获取更多干货