本文提出了uDepot,这是一个KV存储。充分利用了像 Intel Optane这特替快速NM存储设备的结选。杂龄定 研了uDepot结例下它所用的股间中心包备的问题指 并且与现有条统相比可以逐升最利用这些的设备。这外。 本文及表示了uDepot可以使用这些是未来是一份是 服务。很整分的性能与其于ORAA实现的维升变成。但所 条架信息等。表现上,已经使用中心的特别和比如为企 现作为一个实验性公有云服务(39)的基础。

uDepot有两个主要的限制,计划在未来的工作中加以解 以上,当先,也与中时并不《成效地》支持一些已经被证明 对应用程序有用的一些操作。比如范围直接,事务、检查 点、数据结构抽象等(78)。其次,有很多机会可以通过支 持多程户(13)来提高效率,而uDepot目前并没有利用到

贡献

8988

Flash KVZZ68

DRAM KAZERE

NUM KUZYAR

1) 盖人干万(510)个条目,其中没有发生重新调整 大分级件(2) 基人10亿(100)个条目,其中发生了4 有相信服务。进行全等发现,通过信气中的资金等级。 互复(universal benchman),为千万万和10亿条目录 这的股底了20亿元的股份的是,是被逐渐所示,对千万万名。 大约本万万名。大约本级市场下场的大约。 1000年间至2000万万分的

索引结构

嵌入式uDepot

绘出了所有工作负载下的256个客户机块理的高社量。 uDepot使用trtspdk后就提高了YCSB的否社量,相比 Aerospac分据题高195倍(工作负载D)和2.1倍(工作 负载A),相比ScyluDB分别距高10.2倍(工作负载A) 和14.7倍(工作负载B)。

uDepot Server/YCSB

存储设备空间管理 =

Reaping the performance of fast NVM storage with uDepot 借助uDepot获得快速NVM存储的性能

log-structed方式。空间被连续分配。GC处理碎片



1. 在NAND闪存这种特殊存储器上有较好性能。

2. 即使对DRAM这种非特殊存储器,它也比传统方法更高效。

一个重要的用例是镀存。在协同设计GC和镀存时有大量的优化机会。

修改一:使用条目整量的二次高对於新表进行查引。學似于组相联復得[38]:使用从key计算而来的指议的最低有效位 (158) 来计算邻城条引。这允许在重新调整大小明同角效地重度数划路设。而不需要完全存储部设或执行 [0亲获取key并重新计算。

修改二:本文沒有健护每个邻域的位图,也沒有健护每个 邻域的所有条目的继表。而这两者是原始算法所建议使用 的1371。后者将使默认思路的均存需求增加50%(5字节 的条目,都近大为522。每个8日4字节),一个继表至 少需要两倍的均存(保设使用8字节部针和单键表或双键

哈希表 = 表):更不用说增加的重杂性了。本文不使用位图或链表,而是直接对条目的lookup和insert提作执行线性探测。

**南北 =** 

索引数据结构 = 哈希表条目 =

容量利用率 三



uDepot实用C++11实现。值得注意的是,uDepot的性能需要许多优化:使用本地CPU核心的slab分配图从数据路径滑镀金分配,使用大内存页面(huge pages),与动态与相比逻辑的于静态。混合使用scatter型U带来的数据重制,在IO缓冲区的适当位置存放网络数

实现说明

性能/成本比例位于DRAM和SSD之间,平衡了目前流行体系结构将数据全部存放在主存中的趋势。

快NVM设备 (FNDs) FNDs格性能瓶颈从存储转移到网络,因此FNDs KV可以 提供和DRAM一致的性能。

现有KV将数据全部放在DRAM上需要OS显式的以页的方

将数据故在存储设备上的键值存储,即使是那些专门针对 传统SSDs设计的存储器也考虑到不同的需求:较慢的设 备,较小的容量,是否需要多设备或多核扩展。 现有KV store无法充分利用FNDs

一个自腐向上设计的利用FNDs性能的键值存储。uDepot 是一个嵌入式存储可以被应用当作一个库进行使用。

构建了两个网络服务:一个是使用自定义网络协议的分布 式XV存储。另一个是实现Memcache协议的分布式缓 存,可以作为memcached的完全替代,Memcached是 一个广泛使用的基于DRAM的缓存。

1. 实现低延迟

2. 提供每个核心高带宽 3. 在多个核心和设备 F扩展器件部

4. 在字节和操作数上,限制端到端IO放大的较低

5. 存储容量的高利用率

多核优化

背景和动机

存在的问题

有效利用FNDs。许多现存的KV存储使用同步IO严重降低 了性能,因为它依赖于内核调度的线程来处理并发性。 uDepot采用异步IO并在可能的情况下直接从用户态访问 存储设备。为此,uDepot建立在TRT (task runtime)

uDepot使用一个高性旅的DRAM索引性的能够还能FNDs 的性態。同时性特殊内容和饮物、(小内容和饮物效 放的容量利用器。由于较小SDAM需要索引力传统等 量。)如epo的索引性热量可需求小。为了根据存储 家的数目来调整其内存消耗。调整大小不需要任何心操 作,以调量的方式块行使伸手能量小体。

1. uDepot,一种传递了FNDs性能的KV存储,提供低延迟,高吞吐,扩展性和高效的CPU,内存、存储利用率。

2. TRT,一个任务运行时间系统,适合微秒级的IO,作为 uDepot的基石。TRT为写能够开发快存储的应用提供了

 uDepot的索引数据结构能够满足其性能目标。同时在空间上高效。动态可调整使得满足KV对存储的数量。 帝献

4. 实验评估结果证明Depos符合FNDs的性脑,被我们 所知,还没有系统细胞模型,实际上,Upepotit/K化 SSD存储的性能符多,可达14/6. 泛正配基于DRAM 的memsached服务性能,允许将其用FMemsache进行 替缺以重整管研究本,使用uDepo的云Memcache提合 作为实验性理供在公共云上。

DRAM尺寸正在接近物理极限,DRAM也越来越贵。

随着容量需求的增大,内存KV存储依赖模向扩展,通过 增加更多的服务器来实现所需的存储容量。

由于节点的其余部分 (CPU, 存储) 未得到充分的利用。 格外节点也需要按比增加资源如空间,电源,冷却设备。

高效利用FNDs的键值存储提供了具有吸引力的可选项作为DRAM的替代品,特别在普通的网络环境中(10Gbits/s的因特周下)FNDs将瓶颈从存储转移到了网络,网络 环境无法发挥出DRAM键值存储的全部性能。

使用一个多核多设置的系统(20.ore.3.24W/Me drive)与两个单型性的影响的影响的影响的影响的影响的影响的影响。 使用Linux sintsPoxima,到底就是为50M4K的健康的,使用Linux sintsPoxima,到底就是为50M4K的健康的,在对bookstoBHWWedTiger及为1945后,仍无法超过Meprix和20xepss,另一种1945后,但是1945年的1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950年间,1950

1. 基于慢设备构建的系统使用同步IO在微砂级造成很大

2. 使用LSM或B树索引造成大的IO放大。

3. 在DRAM中级存数配需要额外的同步,同时也由于内存 需求限制了扩展性。

4. 提供了许多额外的功能(如事务)对性能造成不好的 影响。

1.同步IO: 网络中需要线程一对一处理并发请求,在请求数大于内核数时,上下文的切换会导致无法充分利用性

2. 异步IO: Linux AIO允许多个IO在一个线程上批量收发。但仅此并不足以充分利用性能。 三种访问存储

用户空间IO:SPDKi量过避免上下文切换。数据负责和 调度开梢来最大化性能。但由于这需要直接访问设备而经 常无法使用,而且许多环境无法支持。

任务被联合调度,而每个任务有自己的栈。TRT生成许多 线程(一个核一个),并在每个自己的栈上执行用户空间 调度器。

调度器和任务之间的切换是轻量的,包括寄存器的保存和

TRT, 任务运行财系统

由于TRT尽量避免内核问题信,所以为同步提供了两个变 练: \*\*intra-and inter-\*\* core、intra-core更有效由于 不需要同比它因此并发访问,只要关键部分不包含切换到 调度器的命令。

TRT为异步IO提供基础结构。一个典型场景:每个网络连接被不同TRT任务服务。为了启用不同IO后端和工具,每个后端来现一个轮询任务未能测事件并通知任务处理。每个核心运行自己的轮询任务来例。