

分布式 VMM 通信架构的研究与原型实现

宋忠雷, 肖利民

(北京航空航天大学计算机学院, 北京 100191)

摘要: 针对分布式虚拟机监控器(DVMM)的通信需求, 研究并实现一套 DVMM 的通信方案, 利用精简可靠数据协议为分布于多台物理主机之上的 VMM 提供可靠、有序、高效的通信服务。通信方案测试表明, 与现行 TCP/IP 协议栈相比, 尽管该方案的带宽并无提升, 但是对通信延迟却减小了 45% 左右, 显示了该方案的可行性和优越性。

关键词: 分布式虚拟机监控器; 通信模块; 通信协议; 精简可靠数据协议

Research and Prototype Implementation of Communication Architecture in Distributed VMM

SONG Zhong-lei, XIAO Li-min

(School of Computer Science and Engineering, Beijing University of Aeronautics and Astronautics, Beijing 100191)

【Abstract】 This paper presents Distributed Virtual Machine Monitor(DVMM) and because of the uniqueness of requirement, a communication solution is adopted and implemented. Reduced Reliable Data Protocol(RRDP) provides a highly efficient, orderly and reliable communication service for DVMM that runs on multiple hosts. Test results in DVMM compared with the current communication protocol stack show that, although the bandwidth is not improved, the latency is reduced by 45%, so it is feasible and advantageous.

【Key words】 Distributed Virtual Machine Monitor(DVMM); communication module; communication protocol; Reduced Reliable Data Protocol(RRDP)

1 概述

随着大规模集成电路的发展, 现代计算机处理能力越来越强大, 这使得现代计算机具有足够的力量来利用虚拟化技术支持多个虚拟机(Virtual Machine, VM), 提供给每个操作系统一套独立的虚拟物理资源。虚拟化技术指的是对计算机资源的抽象化^[1], 为提供更高的资源利用率和灵活性而将物理硬件与操作系统分开的一种技术。目前 X86 架构下的 VMM 主要有 EMC 公司的 VMWare、微软公司的 Virtual Server 和剑桥大学的 Xen^[2], 这种技术由 IBM 公司于 20 世纪 60 年代发明^[3], 在当时用于用户尽可能充分地利用大型机资源。

然而, 近年来, 随着集群的出现以及高性能计算的需要, 一些特殊应用需要更强大的物理资源, 单宿主机已经满足不了应用需求, 将多个主机物理资源有效地进行整合并提供给上层应用成为必然趋势。分布式虚拟机通过分布式虚拟监控器(Distributed Virtual Machine Monitor, DVMM)对硬件资源进行管理, 并通过彼此间的协作提供给客户操作系统统一的物理资源, 从而满足高端服务需求。

DVMM 技术将多台物理机资源进行整合, 为运行在上层的 GOS(Guest Operation System)提供统一的资源, 从而充分利用物理资源, 满足应用需求。Xen 是一个优秀的开源虚拟机, 本文讨论的 DVMM 是基于 Xen 实现的, 它由各物理机上运行的 VMM 组成, 各个 VMM 包括了系统启动模块、指令虚拟模块、I/O 虚拟化模块、中断虚拟化模块、内存虚拟化模块和通信模块。各个模块除了独立对硬件资源进行虚拟化外, 还需要通过通信连接起来。VMM 之间同样需要通信进行连接整合, 为 GOS 呈现出一台“物理主机”资源, 如图 1

所示。

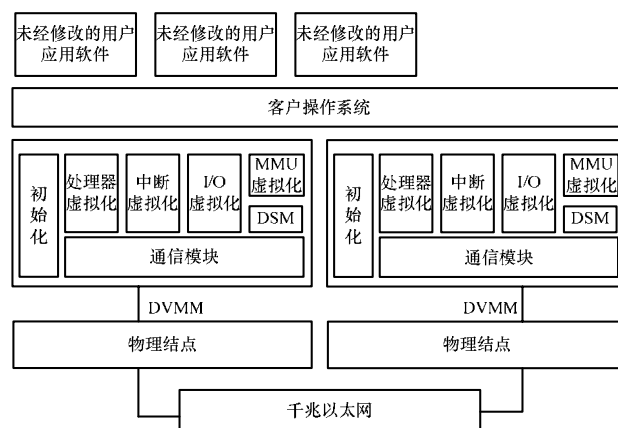


图 1 DVMM 各模块视图及通信系统所处的位置

实现 DVMM 的关键之一是为运行在不同主机上的 VMM 之间提供高效可靠的通信。由图 1 可知, 通信模块是分布式 VMM 的基础, 它为不同 VMM 的各个模块之间提供了协作的平台。各物理主机的 VMM 通过底层通信软件传递控制信息和数据信息, 从而完成分布式 VMM 的各项功能, 如: 为操作系统实现单一内存空间以及跨节点的处理程序统一调度等。

基金项目: 国家“863”计划基金资助项目“新型服务器体系结构及其关键技术”(2006AA01Z108)

作者简介: 宋忠雷(1984-), 男, 硕士研究生, 主研方向: 计算机系统结构, 虚拟化技术, 网络通信; 肖利民, 教授、博士生导师

收稿日期: 2009-11-30 E-mail: szl_84@sina.com

- 2)ACK, 确认域有效
- 3)EACK, 乱序确认
- 4)RST, 重新开始建立链接
- 5)NUL, 表明该帧没有有效数字
- 6)第 5 位没有使用
- 7)第 6 位和第 7 位是协议的版本号
- (2)Header Length, 协议头部的长度
- (3)Source and Destination Ports, 源端口和目的端口号
- (4)Data Length, 有效数据的长度
- (5)Sequence Number, 包的序号
- (6)Acknowledgement Number, 确认序号
- (7)Checksum, 校验和

在以上各字段的共同作用下, RRDP 完成连接建立与关闭, 以及正确传输数据的功能。

	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
0	SYN	ACK	EACK	RST	NUL	0	VerNo.	Header Length								
1	Source Port															
2	Destination Port															
3	Data Length															
4	Sequence Number															
5	Acknowledgement Number															
6	Checksum															

图 4 可靠传输 RRDP 协议首部

3.3 协议实现

根据上文设计, 本文实现了 RRDP 协议, 并加入了 DVMM 系统中, 使其可以为 DVMM 系统提供可靠、有序、高效的通信服务。

图 5 为 RRDP 协议实现层次视图, 与 TCP 和 UDP 协议栈相比, RRDP 协议更加简洁, 由于是局域网内的协议, 去除了 IP 层, 体现了其专用性。而且, 在协议实现过程中, 为提高数据传输的效率, 使用精简的数据结构和优化的算法。在通信协议中, 无论是发送数据还是接收数据的过程都会涉及到很多通信数据结构和通信控制结构, 通过精简这些结构和优化处理这些结构的算法能适当提高协议的传输效率。

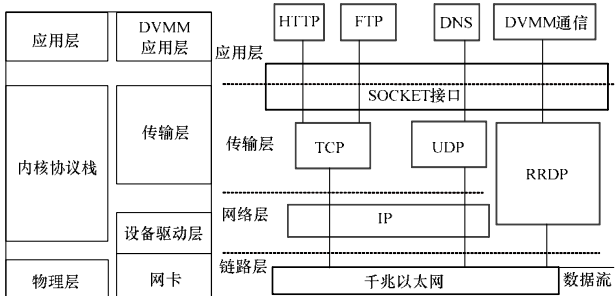


图 5 RRDP 协议层次及实现结构

其次, 减少数据在主机中的拷贝次数, 对于 TCP 或者 UDP, 数据首先从用户态拷贝到内核态, 在 IP 层对数据块进行分片时会再一次进行数据拷贝。而在 RRDP 协议的实现中, 通过用户态与内核态共享数据的方法实现用户到内核的零拷贝, 在去除了 IP 层之后, 传输层通过 DMA 技术将数据直接交付给网卡处理从而达到了真正的零拷贝。这样极大地提高了协议性能。

3.4 地址转换层

为保证与标准 socket 接口兼容, RRDP 模块必须提供简单的地址转换协议, 它负责将 DVMM 系统中各个物理结点的结点号与它们的 MAC 地址对应起来。转换模块负责维护 MAC 表(保存物理结点号与 MAC 地址之间的对应关系), 在通信的建立连接过程中, 各个 VMM 使用它们各自的结点号, 这样提供了高效的地址转换方式。

4 性能测试与结果分析

本文针对 DVMM 系统通信应用进行了性能测试, 环境为 2 个具有千兆以太网的服务器运行 DVMM 系统以及 Dom0。dm 为设备管理模块, 协议以模块形式加载到 Dom0 中, 测试程序模仿 DVMM 应用环境编程。测试模型如图 6 所示。测试内容主要包括点到点的单向通信延迟和通信带宽, 并与 TCP/IP 协议栈的相应性能进行了比较。

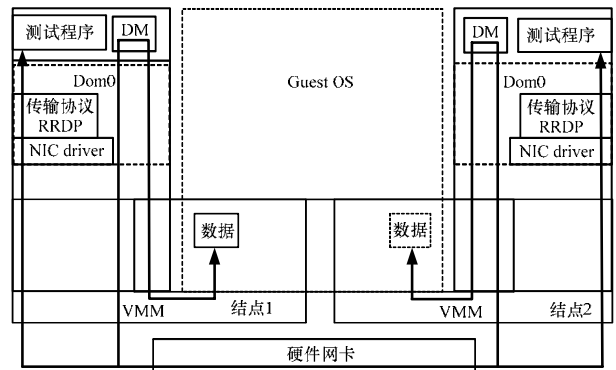


图 6 通信测试模型

4.1 性能测试

4.1.1 通信延迟测试

图 7 显示了 DVMM 通信协议 RRDP 延迟和传输数据块大小的关系, 并与 TCP 延迟进行比较。

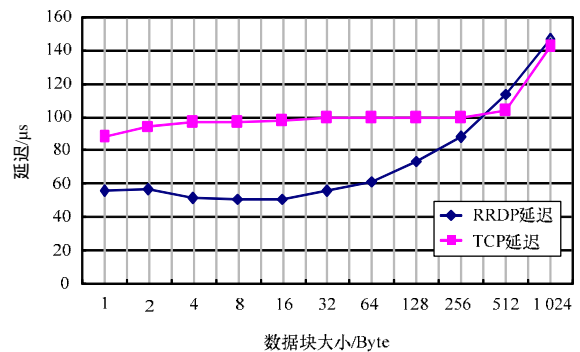


图 7 RRDP 与 TCP 协议延迟比较

从图 7 可以看出, 随着数据块不断增大, 通信协议传输延迟逐渐增大。而对于小数据块(32 Byte 以下)增加不太显著, 这主要是因为网络传输的小数据帧中头部占得较大(RRDP 头部和 MAC 头部), 而随着数据所占部分增大延迟明显大幅增加。与 TCP 协议相比, RRDP 协议延迟大概减小 45% 左右。当数据块大于 512 Byte 时, RRDP 较 TCP 协议延迟减小不太明显。因为数据帧较大时, TCP 和 RRDP 的延迟大部分都消耗在了用户层到内核层之间的数据拷贝, 使得两者差别不大。

4.1.2 通信带宽测试

图 8 给出 DVMM 通信协议 RRDP 带宽和传输数据块大小的关系, 并与 TCP 进行了对比。由于 RRDP 无分片功能, 因此在千兆网卡大帧支持下最大数据块为 9 000 Byte。

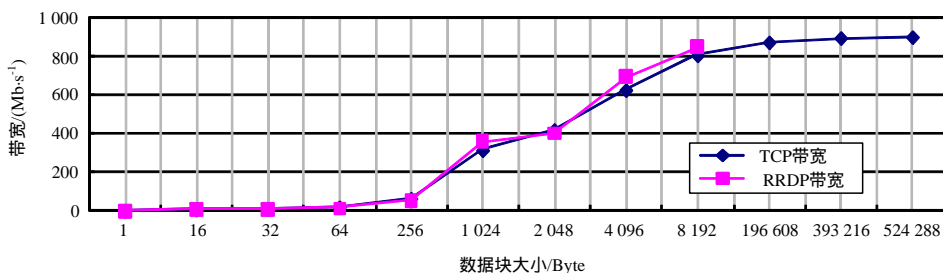


图8 RRDP与TCP协议带宽比较

在通信带宽测试中,同一环境下,RRDP和TCP协议带宽均逼近于900 Mb/s,两者差异不太明显。在DVMM通信中,提高通信协议质量关键在于提高延迟,因此,对于带宽,意义并不是太大。

4.2 结果分析

分析可知,实验中在数据字节数较小时延迟有很大提高,这对DVMM系统有很大意义。因为在DVMM系统通信中,多数的信息交换都是小数据块传输,如中断的发送、BIOS信息交换、远程I/O请求等。

而对于大数据块可以通过零拷贝技术进一步减小延迟。虽然带宽无显著增大,但通信方案最主要的目的——低延迟已经达到。作为简洁的三层协议栈,具有一定的潜力,但是通用性不强,可考虑在加入网络层的基础上使其成为与TCP、UDP对等的因特网协议。

5 相关工作

5.1 MM间的通信

日本东京大学的Virtual Multiprocessor项目也是实现分布式虚拟机。在这个项目中,在VMM和硬件层之间,运行了一个简化的操作系统,称为Host OS。VMM间的通信使用自己的应用协议,它调用Host OS的TCP/IP协议栈来完成数据的发送与接收。

目前常见的局域网通信协议主要有NetBEUI、IPX/SPX、NWLink、TCP/IP。在这几种协议中用得最多、最为复杂的是TCP/IP协议。这些协议的实现都离不开操作系统的支持。以最常见TCP/IP协议为例,其实现复杂,不仅可以实现局域网通信,而且支持“路由”功能。因此,不适合VMM对通信的要求。

(上接第112页)

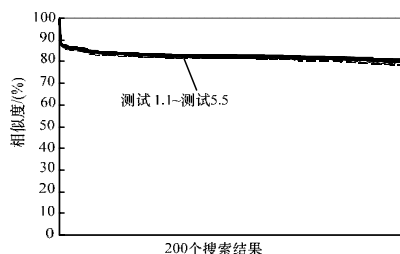


图5 扰动环境下SWAD的仿真结果

6 结束语

本文对SWIM算法进行了抗扰动改进,实现了基于多媒体特征的抗扰动P2P搜索。在发布文件并与邻居文件交换更新信息的过程中,如何减少重复数据的发送,提高通信效率,仍有很大的改进空间。

5.2 其他几种典型的通信

VMMC是一种基于虚拟内存映射的通信机制,它支持从发送方虚拟内存到接收方虚拟内存的数据直接传送。

Active Message是一种异步通信机制,采用与传统的通信机制完全不同的设计思路,更直接地使用通信硬件提供的功能。

Fast Sockets是一种面向局域网的通信软件,它采用新的高效通信协议在用户空间实现了一般Unix操作系统均提供的Berkeley Sockets API。该软件试图通过修改现有的网络通信协议,使得通信过程中软件开销最小,同时又能保持Fast Sockets与现有应用程序以及广域网通信协议的兼容性。

6 结束语

本文提出并实现了分布式VMM间的通信架构,为实现分布式VMM提供了基础。它的性能是影响分布式VMM性能的关键因素之一,因此,研究如何进一步提高底层通信的性能就显得十分重要^[4]。通过分析得知,由用户态到内核态的数据拷贝影响了通信的性能,后续将实现通信部分零拷贝技术来提高它的性能,采用Infiniband高速网代替千兆以太网等策略来进一步优化分布式VMM底层通信的性能。协议直接运行在网卡驱动层之上,通过修改驱动也能从一定程度上提高它的性能。

参考文献

- [1] 董耀祖,周正伟. 基于X86架构的系统虚拟机技术与应用[J]. 计算机工程, 2006, 32(13): 71-73.
- [2] Barham P, Dragovic B, Fraser K, et al. Xen and the Art of Virtualization[C]//Proc. of the 19th ACM Symposium on Operating Systems Principles. Bolton Landing, NY, USA: Virtual Machine Monitors, 2003: 164-177.
- [3] Creasy R J. The Origin of the VM/370 Time-sharing System[J]. IBM Journal of Research and Development, 1981, 25(5): 483-490.
- [4] 付赛平,任国林. XEN网络I/O完全虚拟化机制的可扩展性研究[J]. 计算机工程, 2008, 34(23): 102-104.

编辑 顾逸斐

参考文献

- [1] Novak D, Zezula P. M-Chord: A Scalable Distributed Similarity Search Structure[C]//Proceedings of INFOSCALE'06. Hong Kong, China: [s. n.], 2006.
- [2] Androutsos P, Androutsos D, Venetsanopoulos A N. A Distributed Fault-tolerant MPEG-7 Retrieval Scheme Based on Small World Theory[J]. IEEE Transactions on Multimedia, 2006, 18(2): 278-288.
- [3] Maymounkov P, Mazières D. Kademia: A Peer-to-Peer Information System Based on the XOR Metric[C]//Proceedings of the 1st International Workshop on Peer-to-Peer Systems. London, UK: [s. n.], 2002.
- [4] ISO/IEC. ISO/IEC 15938-6:2003 Information Technology — Multimedia Content Description Interface-Part 6: Reference Software[S]. 2003.

编辑 索书志