

一种高效 GPU 存储系统体系架构设计

卢俊, 颜哲, 田泽

(中国航空计算技术研究所, 陕西 西安 710068)

摘要: 图形处理技术被广泛应用于电影、视频、游戏以及动画的制作, 而图形处理系统 (GPU) 的出现极大地减轻了 CPU 日益繁重的图形处理任务, 使得其能更专注于通用控制。文中阐述了制约 GPU 性能提升的重要因素, 指出提高带宽利用率是应对这一问题的关键措施。通过局部性原理的分析, 提出了一种基于层次化架构的高效 GPU 存储系统的设计。文中介绍了 4 层结构的存储系统, 并逐层说明了各自的功能和架构, 评估了基于层次化存储架构的 GPU 在典型应用中的带宽。文中还描述了 Cache 以及显存管理等子模块的功能。通过仿真可知, 该 GPU 存储系统能充分利用共享和复用等手段尽量减少外部存储器的访问次数, 从而提高了带宽利用率。

关键词: 图形处理系统; 层次化存储; 带宽; 存储管理模块

中图分类号: TP31

文献标识码: A

文章编号: 1673-629X(2015)04-0006-04

doi: 10.3969/j.issn.1673-629X.2015.04.002

An Efficient Memory System Structure Design of GPU

LU Jun, YAN Zhe, TIAN Ze

(Aeronautical Computing Technique Research Institute, Xi'an 710068, China)

Abstract: Graphic processing technique has been widely used in movie, video, game and cartoon making. GPU has greatly reduced the pressure of graphic processing, which was CPU's job in the past. It introduces the key factor which limits improvement of GPU performance, and also indicates that the bandwidth utilization is one of the critical resources to deal with that limits. According to analysis of principle of locality, attempt to explore an efficient GPU memory system based on layered architecture. This global memory hierarchy structure has four layers of memory, and the function of each layer has been described. In this paper, present the bandwidth utilization ratio of typical GPU running scenarios, also introduce the architecture of the modules such as Cache and MMU. This GPU memory system can reduce the frequency of accessing SDRAM by sharing cache on chip, so that the utilization ratio of bandwidth has been greatly improved.

Key words: graphic processing system; memory hierarchy; bandwidth; MMU

0 引言

图形处理系统 (GPU) 的出现极大地减轻了 CPU 日益繁重的图形处理任务, 使其能更专注于通用控制^[1]。然而相比较 CPU, GPU 拥有更多的运算内核、更多的 ALU, 需要存储通路提供更大的带宽, 这使得 GPU 的整体性能严重依赖于存储系统, 尤其是片外存储的数据传输能力。统一染色架构的多核技术在 3D 处理中的应用使 GPU 处理能力得到极大的提升, 然而性能的增长速度远大于存储器存取速度的增加, 导致了“存储墙”问题变得更加严重, 也对存储系统提出了更高的要求^[2]。目前, GPU 与片外存储器之间的带宽需求已经达到 200 GB/s, 存储带宽已经成为进一步提

高 GPU 性能的瓶颈。在有限的硬件资源下, 应对这一瓶颈的有效手段就是提高带宽利用率, 即: 充分使用共享和复用等手段尽量减少从外部存储器中取数的次数^[3-4]。

GPU 对存储器中的某个数据进行访存操作后, 该数据很可能在很短的时间内被再次访问而产生时间局部性; 对存储器中的某个数据进行访存操作之后, 该数据存储地址的相邻地址中的数据也可能被访问而出现空间局部性。通过充分捕捉指令、图像和纹理等数据的局部性, 采用多层级的存储资源可以有效减少对片外存储器的依赖, 达到合理高效地利用存储带宽的目的^[5-6]。

收稿日期: 2014-06-13

修回日期: 2014-09-17

网络出版时间: 2015-03-31

基金项目: 国家“十二五”微电子预研基金项目 (51308010601, 51308010710, 51308010711)

作者简介: 卢俊 (1981-), 男, 工程师, 研究方向为集成电路设计与验证; 田泽, 博士, 研究员, 中国航空工业集团首席技术专家, 研究方向为 SoC 设计、嵌入式系统设计与 VLSI 设计等。

网络出版地址: <http://www.cnki.net/kcms/detail/61.1450.TP.20150331.0941.002.html>

1 层次化存储架构

为了减少 GPU 各请求源对外部存储器的访问频率,减少 DDR3 存储通路的带宽压力,采用层次化的 4 级存储结构,存储架构图如图 1 所示^[7]。L₀层是各请求源内部的寄存器级缓存,速度最快、容量最小、访问最频繁;L₁层主要包括 L₁ Cache 以及视频通路的读写

行缓冲存储器;L₂层是染色程序存储器、纹理 Cache 的二级缓存;L₃层是两条 64 bits 数据位宽的独立通路,由仲裁控制、DDR3 控制器、PHY 以及 SDRAM 存储器芯片组成,用以访问 GPU 的片外存储器,其访问速度最慢、容量最大^[8]。

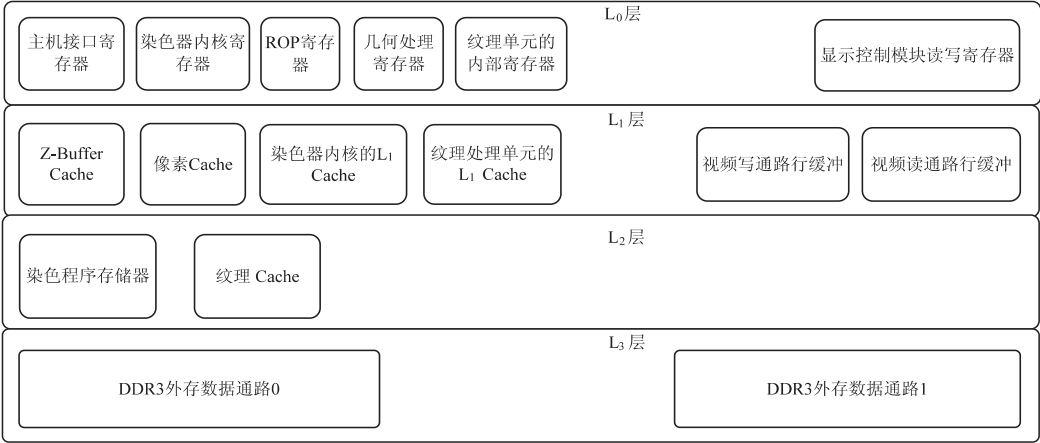


图 1 GPU 的 4 级存储架构

1.1 Cache 的设计

Cache 位于存储系统的 L₀层和 L₁层,包括有四种不同功能的 Cache:Z-Buffer Cache、像素 Cache、染色器内核 Cache、纹理 Cache。用于实现 GPU 运算单元寄存器层与片外存储器之间的数据读写缓冲^[9-11]。

当 Cache 接收到 GPU 的读写请求时,执行如下共同的操作:

- (1)判断上述各请求源是否地址命中;
- (2)如果命中,则直接对 Cache 中的 RAM 进行读写操作;
- (3)如果不命中,则首先检测 RAM 的脏位标志,标志为 1 表示有脏位,则将当前块数据回写到 GPU 主存储中,然后替换当前块数据,最后直接对 Cache 中的 RAM 进行读写操作;脏位标志为 0 表示没有脏位,直接替换当前块数据,并对 Cache 中的 RAM 进行读写。

Cache 的主要工作流程如图 2 所示。

1.2 视频读写通路行缓冲

工作流程描述如下:

(1)写操作通路采用双缓机制,将视频源数据按分辨率大小逐行写入缓冲区。当一帧图像的某行数据写完成,向外存数据通路仲裁器发出写请求信号,如果此时 DDR3 空闲则直接响应请求。由于视频源数据按二维地址映射方式存储,因此一行的图像数据能很快写入 SDRAM 芯片。当写请求信号发出时 DDR3 繁忙,则本次写请求进入仲裁队列,此时下一行视频数据被写入第二个缓冲区。由于 DDR3 的 burst 写操作速度远大于视频源写入缓冲区速度,因此双缓机制能保证

视频源的连续写入。

(2)读操作通路共有 8 个缓冲区,正常读请求被响应时会从 8 个缓冲区中依次读出请求的行号对应的数据。当图像需要旋转 90°时,根据显控模块提供的行列信息从 DDR3 SDRAM 中预读取出 8 列图像信息,然后再根据图像向左、向右的翻转类型拼出行数据,此时显控模块所需的视频数据会从缓冲区中被快速读出。当 8 个缓冲区的数据均被读出后,重新预取 SDRAM 中的视频数据。

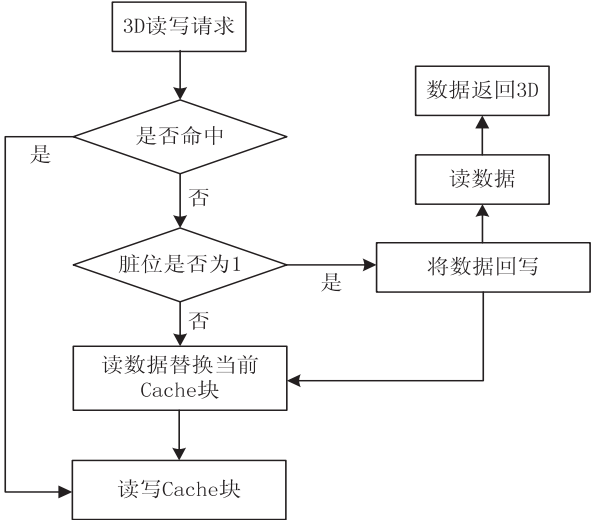


图 2 Cache 工作流程

1.3 显存管理模块的设计

显存管理模块位于 L₃层,由视频通路、图像通路、仲裁控制、DDR3 控制器以及 SDRAM 存储器芯片等模块组成,用以访问 GPU 的片外存储资源。与 Cache 相

比,其访问速度慢、存储容量大。显存管理模块 (MMU)的体系架构如图 3 所示。

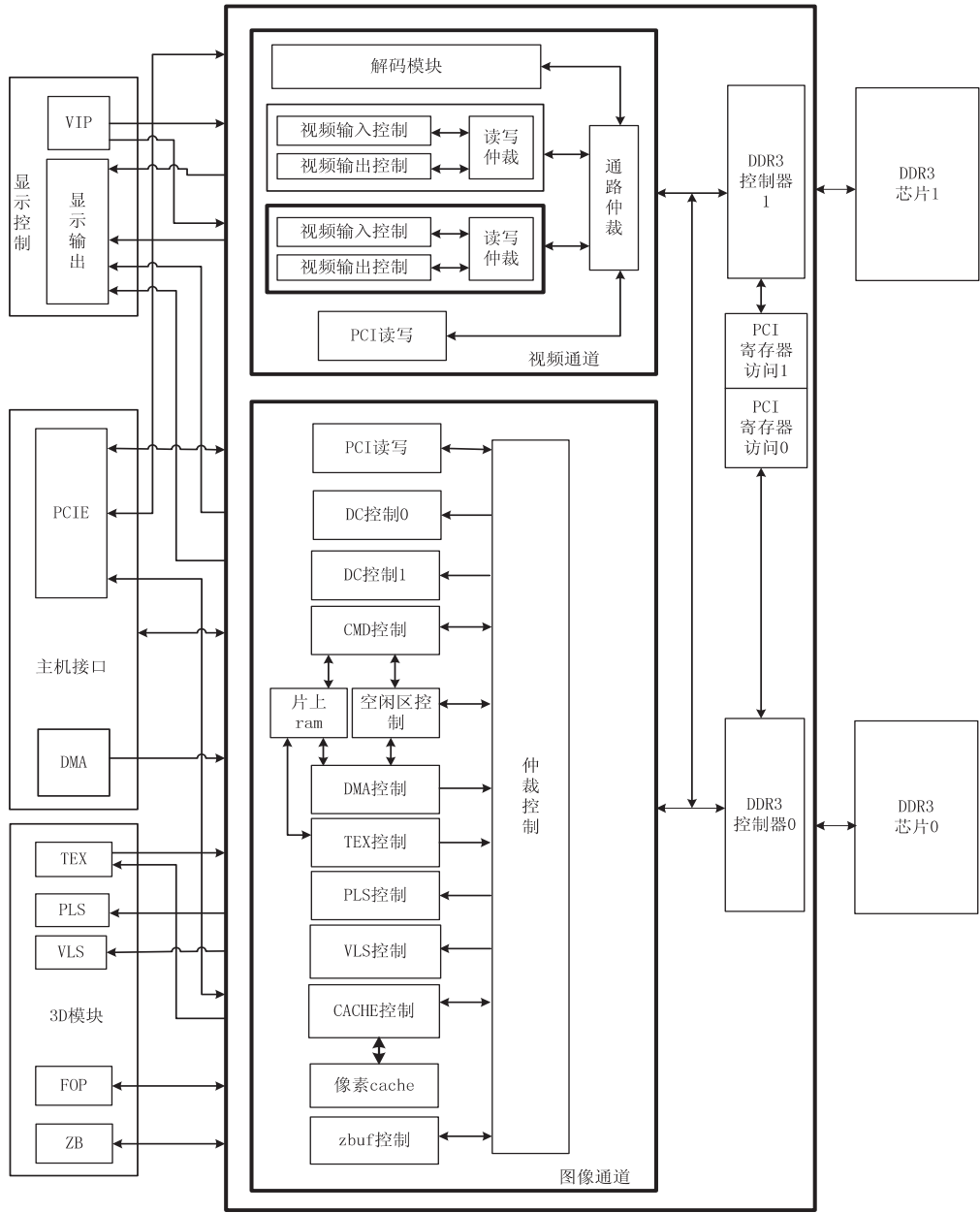


图 3 显存管理模块架构图

GPU 的外部存储管理模块使用两路 64 bits DDR3 SDRAM 作为系统的主存,用于图像数据和视频数据的缓存。访问图像存储管理通路的功能单元主要有命令处理器 DMA、3D 绘图引擎和显示控制模块。DMA 模块主要负责将染色器固定程序和用户程序写入存储器;3D 绘图引擎负责在绘图过程中从存储器读取数据、存入中间图形数据及绘制完成的图形数据。该通路实现对像素数据的缓存、对各外围单元对 SDRAM 访问的仲裁和存储控制等。图像存储管理通路通过 4 kB 的 Cache 和存储仲裁完成系统中 3D 引擎、显示控制 DC、DMA 等功能模块对 DDR3 的有序访问,使图形图像得以流畅、完整的显示,提高整个流水线的速度,最大化处理器绘图性能^[12-13]。

在视频数据缓存的通路中,访问存储器的功能单元主要有显示控制模块和 H264 解码模块。外存控制模块负责将视频源数据存入 DDR3 存储器及读出并显示输出视频数据;H264 解码模块负责控制解码数据,并将解码后的视频数据存储在 DDR3 存储器中。

存储管理模块的两路 DDR3 控制器架构、功能和特性完全相同,所支持的 DDR3 SDRAM 芯片大小可配置,单路容量最小为 1 GB,主要完成对片外 DDR3 SDRAM 存储器的初始化操作和读写访问控制。

2 带宽评估

采用层次化存储系统评估数据带宽,典型应用的场景如下:3D 引擎每秒绘制 30 帧,按照每像素点 32

bit,场景复杂度按 5(每帧中的每个像素渲染 5 次)计算,每帧分辨率按 1 920×1 440 计算。所需的带宽如表 1 所示。

表 1 GPU 存储带宽需求

序号	请求源描述	特征描述	带宽需求 (GB/s)
1	纹理 Cache	绘图使用纹理时,染色器需要访问存储	0.94
2	Z-buffer Cache	所有绘制的像素点都需要访问存储	0.332
3	ROP 阵列	累积、混合、逻辑操作需要访问存储	1.011+0.8 =1.811
4	像素 Cache	所有绘制完成的图像都需要通过像素 Cache 写回帧缓存	1.659
5	顶点数组和显示列表	常用、数据量较小	0.0231 6
6	显示控制模块	实时显示需要持续访问存储	2.656
7	H264 解码模块	高清 25 帧时的带宽要求	1

运算过程中 GPU 所需的数据带宽为 8.43 GB/s,1 路 64 bits DDR3-1333 类型的 SDRAM 芯片作为外存,带宽利用率按 50% 计算,能提供约 5.2 GB/s 的带宽。为满足 GPU 的带宽需求应采用 2 路外部存储通道(能提供 10.4 GB/s 带宽)。采用上述 4 级存储结构的层次化存储系统能满足在典型应用的带宽需求。

3 结束语

综上,可以看到寄存器文件、片上一级 Cache、片上二级 Cache 以及显存管理模块共同组成了 GPU 存储架构的基本层次。在多级存储层次中^[14],寄存器文件作为离 GPU 数据通路最近的存储器,对处理器综合性能的影响最明显,因此在允许的硬件预算下增大寄存器文件可以有效增强数据局部性和访问速度。片上缓存的作用是填补 GPU 与主存之间在速度上的差距,其本质上是存储器中经常被 GPU 访问数据的一个副本,但缓存需要解决映像规则、查找方法、替换算法、写策略等问题,另外缓存的一致性策略通常比较复杂,占用的硬件开销也比较大。片外存储器作为存储体系最底层,通常容量大、速度较慢,用作 GPU 的外部主存储器。层次化的存储架构缓解了寄存器文件与片外存储器之间的速度差距,增强了存储体系的带宽利用率^[15]。

参考文献:

[1] 蔡士杰,宋继强,蔡敏. 计算机图形学[M]. 第 3 版. 北京:电子工业出版社,2007:10-21.

[2] 徐新海,林宇裴,易伟. CPU-GPGPU 异构体系结构相关技术综述[J]. 计算机工程与科学,2009,31(A1):24-26.

[3] Wolf W. High performance embedded computing architectures, applications, and methodologies[M]. New York:Elsevier, 2007.

[4] Yoo Hoi-Jun, Woo Jeong-Ho. Mobile 3D graphics SoC from algorithm to chip[M]. Republic of Korea:John Wiley & Sons (Asia) Pte Ltd,2009.

[5] 马安国,成玉,唐遇星,等. GPU 异构系统中的存储层次和负载均衡策略研究[J]. 国防科技大学学报,2009,31(5):38-43.

[6] 王鹏,伊鹏,金德鹏,等. 基于三级存储阵列缓存高速数据包及性能分析[J]. 软件学报,2005,16(12):2181-2189.

[7] Lindholm E, Nickolls J, Oberman S, et al. NVIDIA Tesla: a unified graphics and computing architecture[J]. IEEE Micro, 2008,28(2):39-55.

[8] Martin M. Token coherence[D]. Wisconsin: University of Wisconsin-Madison,2003.

[9] Johansson M. General purpose computing on graphics processing units using OpenCL[D]. Sweden: Chalmers University of Technology,2010.

[10] Woo R, Choi S, Sohn Ju-Ho, et al. A low-power 3D rendering engine with two texture units and 29Mb embedded DRAM for 3D multimedia terminals[J]. IEEE Journal of Solid-state Circuits,2004,39(7):1101-1109.

[11] 王钰. 多机可缩放性高速缓冲存储器一致性协议分析[J]. 计算机技术与发展,2009,19(2):94-97.

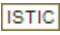
[12] Elder G. ATI Radeon 9700: architecture and 3D performance [C]//Proc of ACM SIGGRAPH/Eurographics. [s. l.]: ACM,2002:86-92.

[13] Garcia J, March M, Cerda L, et al. On the design of hybrid DRAM/SRAM memory schemes for fast packet buffers[C]//Proc of HPSR. [s. l.]: IEEE Computer Society,2004:15-19.

[14] 田绪红,陈茂资,田金梅. DirectX 发展及相关 GPU 通用计算技术综述[J]. 计算机工程与设计,2009,30(23):5432-5436.

[15] 吴俊杰. 层次存储的访问分析与优化方法研究—重用性、相似性与亲和性[D]. 长沙:国防科技大学出版社,2009.

一种高效GPU存储系统体系架构设计

作者：[卢俊](#)，[颜哲](#)，[田泽](#)，[LU Jun](#)，[YAN Zhe](#)，[TIAN Ze](#)
作者单位：[中国航空计算技术研究所, 陕西 西安, 710068](#)
刊名：[计算机技术与发展](#)
英文刊名：[Computer Technology and Development](#)
年，卷(期)：2015(4)

引用本文格式：[卢俊](#).[颜哲](#).[田泽](#).[LU Jun](#).[YAN Zhe](#).[TIAN Ze](#) 一种高效GPU存储系统体系架构设计[期刊论文]-[计算机技术与发展](#) 2015(4)