

# GPU Computing Technology Overview

赵开勇

香港浸会大学/浪潮GPU高性能开发顾问

CUDA社区版主

<http://blog.csdn.net/openhero>

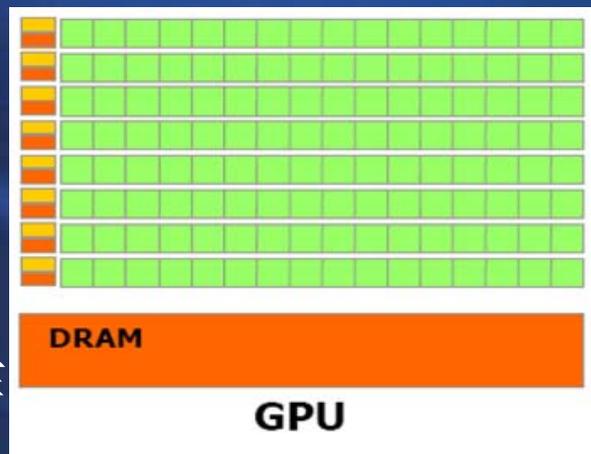
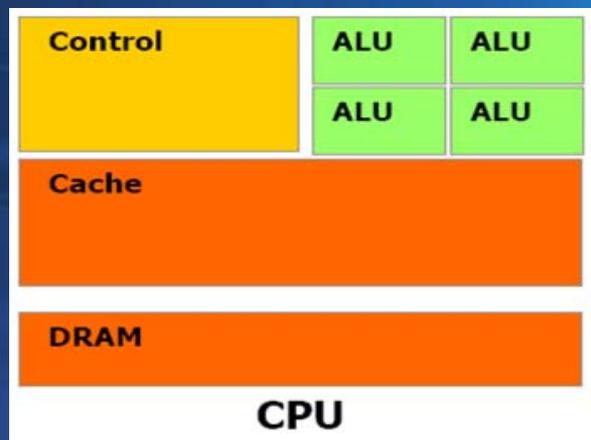
# 提纲

- 处理器核心技术趋势和现状
- 硬件架构简介
  - Nvidia: G80, GTX200
  - AMD(ATI): R600, RV770
  - Intel(Larrabee)
- 软件编程模型简介
  - Nvidia: CUDA
  - AMD(ATI): Stream computing
  - Intel: Ct
- 应用与优化

# 处理器核心技术趋势和现状

## ● 处理器核心技术趋势

- 高效的计算能力
  - 利用数据的并行性
- 高效的数据带宽
  - 增加片内通信，减少片外通信
- 降低指令级的控制复杂度
  - out-of-order vs in-order
- 流水线式处理模式
  - 每一个环节只处理大量简单的工作
- 低功耗



# GPU技术的现状及优势

- 强大的并行处理能力 GPU接近
  - Single float: 1T
  - Double float: 100G~200G
- 高带宽
  - 141GB/s (GTX280), 115 GB/s(RV770)
- 低成本
  - Gflop/\$和Gflops/w高于CPU
- 扩展性
  - CPU的协处理器，现有的PCIE都可扩展

# 提纲

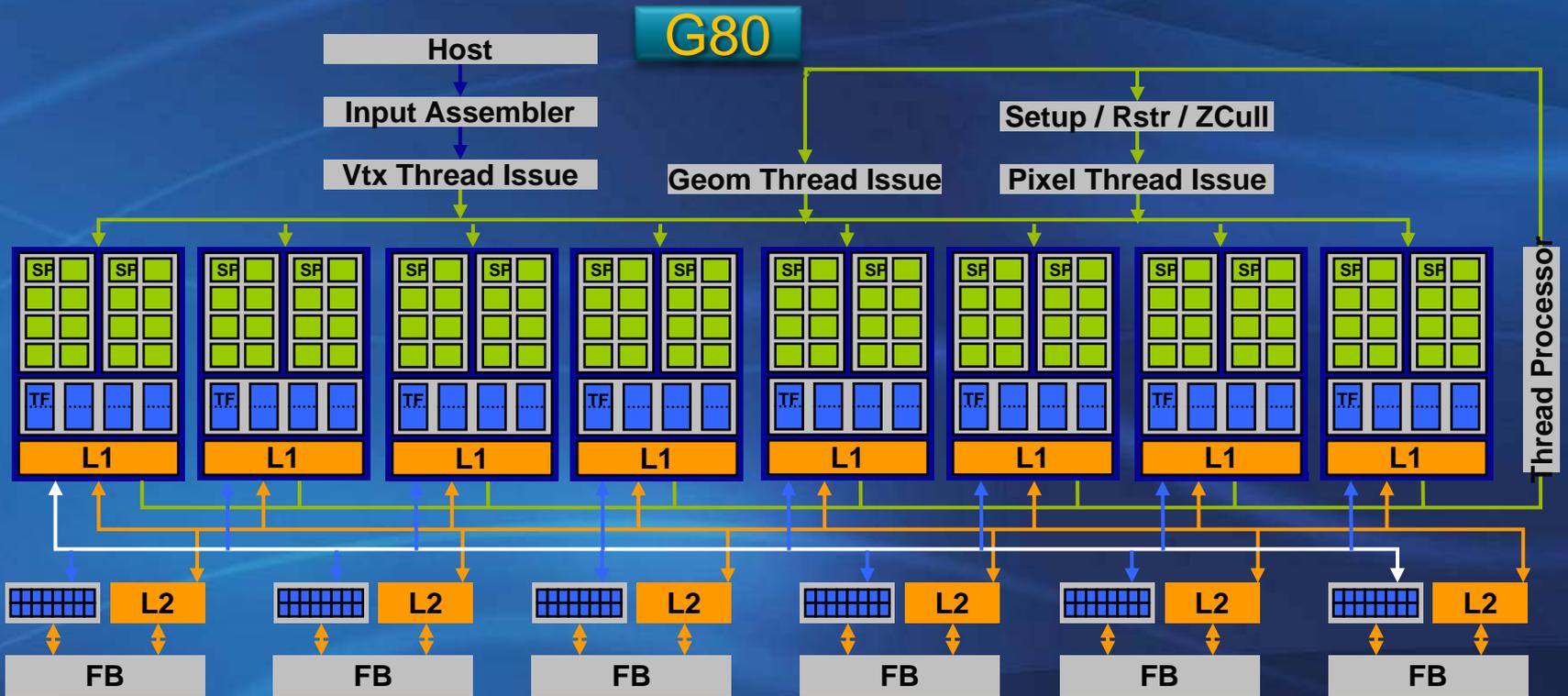
- 处理器核心技术趋势和现状
- 硬件架构简介
  - Nvidia: G80, GTX200
  - AMD(ATI): R600, RV770
  - Intel(Larrabee)
- 软件编程模型简介
  - Nvidia: CUDA
  - AMD(ATI): Stream computing
  - Intel: Ct
- 应用与优化

# GPU技术路线的分析

硬件:

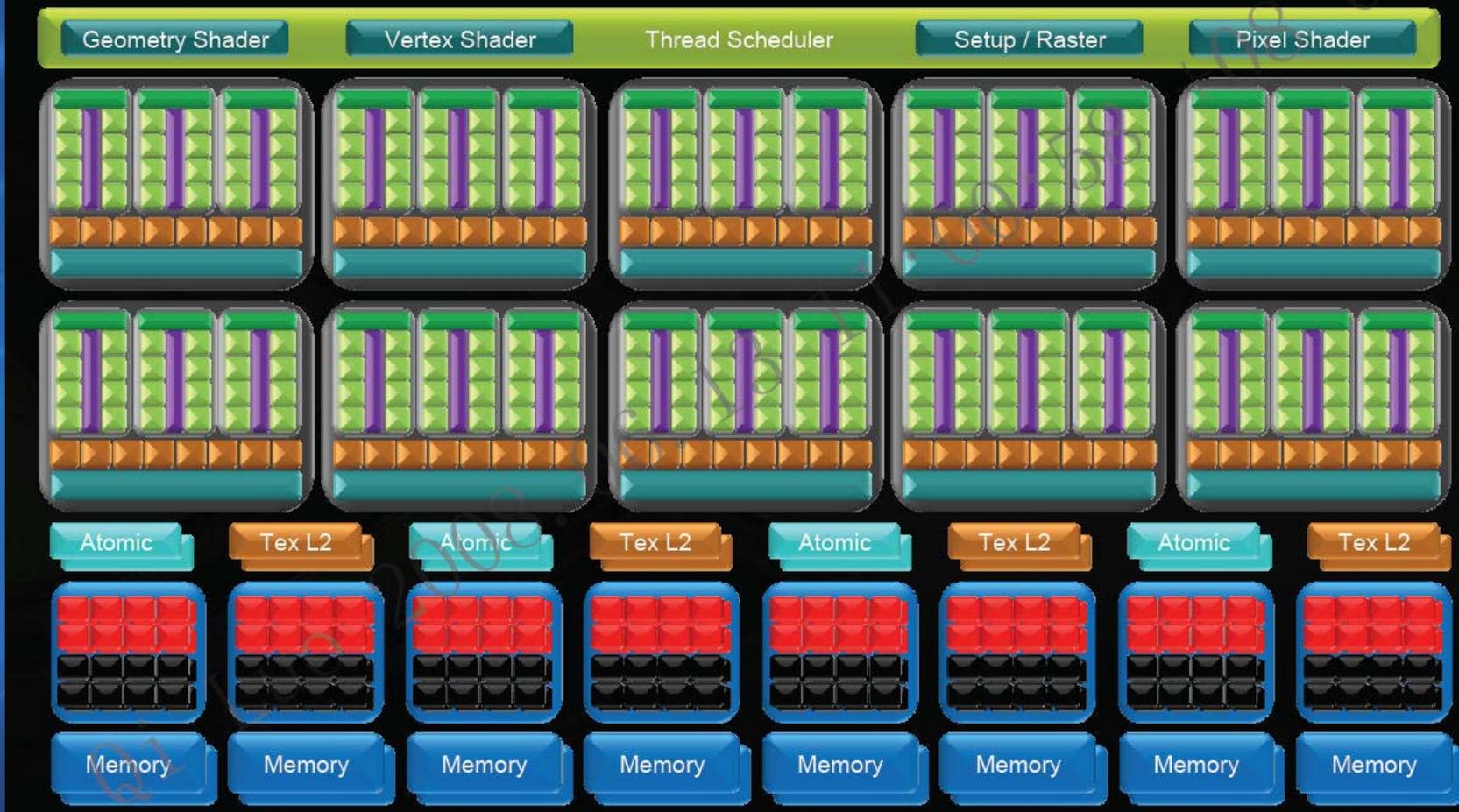
硬件	Nvidia(GT200)	AMD(RV770)	Intel(Larrabee)
单个处理单元	8SP+2SFU+1DU	16x(4ALU+1SFU+1DU)	2 in-order core
总处理单元	240SP	160x5 = 800	32
线程	Thread:SP 1:1	Thread:ALU 1:5	Thread:core 4:2
时钟频率	1296MHz	750MHz	未知
DRAM	GDDR3	GDDR5	GDDR5
内存访问模式	Global+shared	Global+shared	Shared
内存位宽	512bit	512bit	1024-bit

# Nvidia G80、GTX200硬件架构

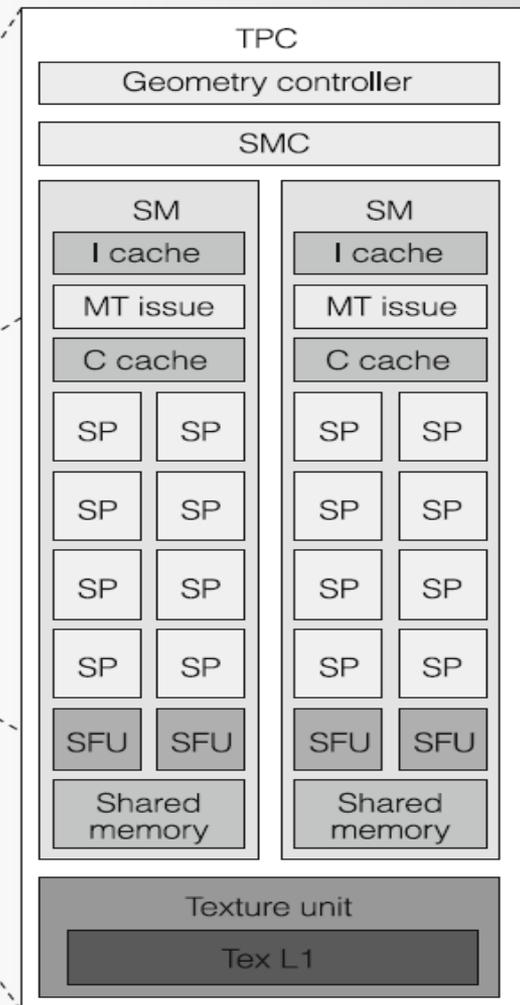
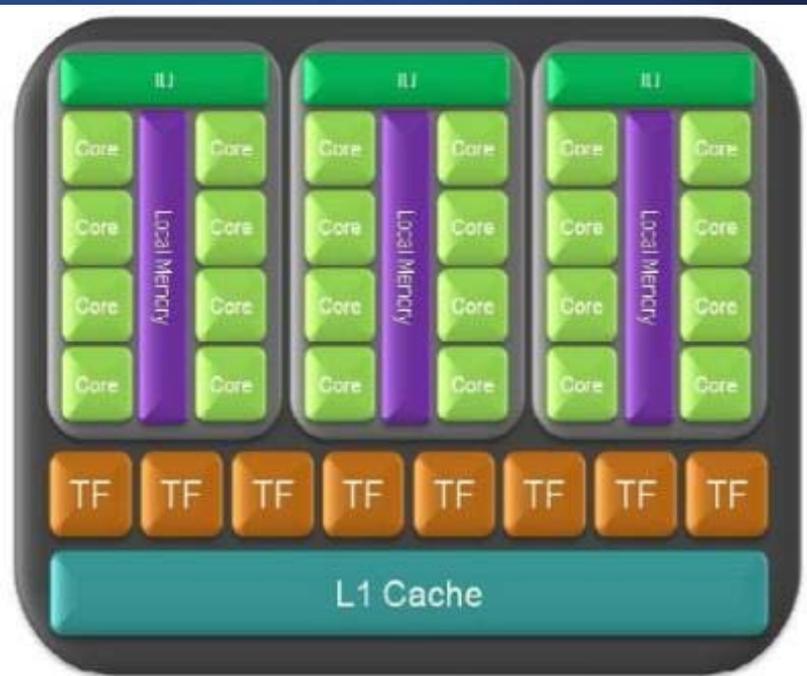


# Nvidia G80、GTX200 硬件架构

## GeForce GTX 200 Architecture



# Nvidia 计算单元(TPC, SM, SP)

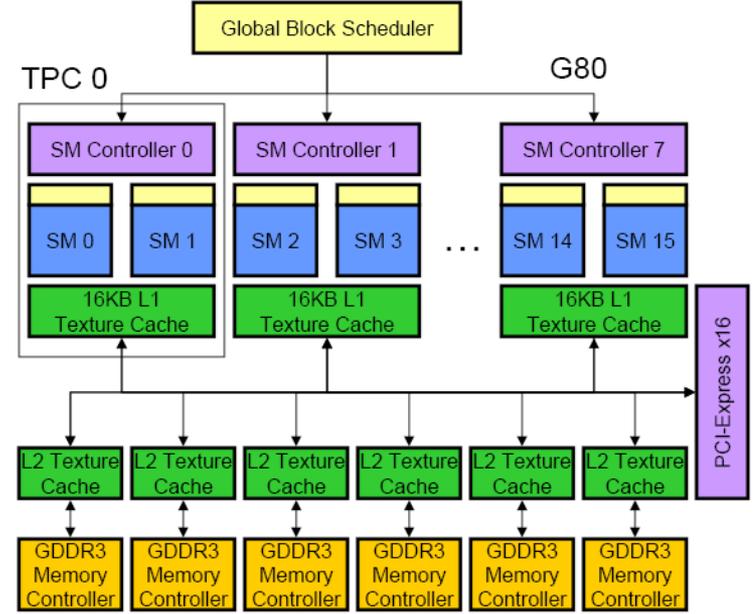
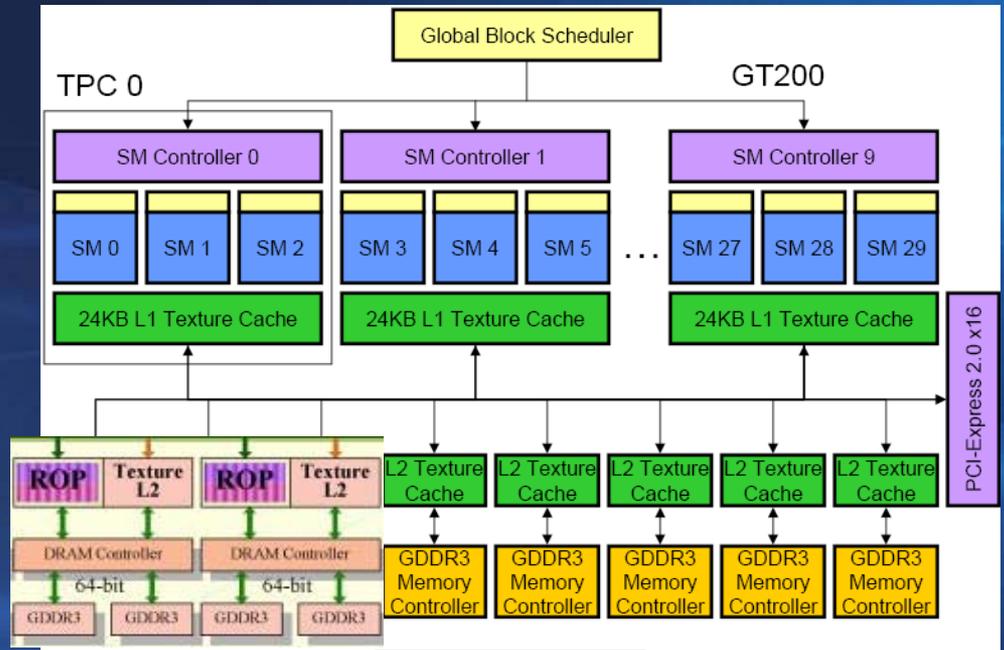
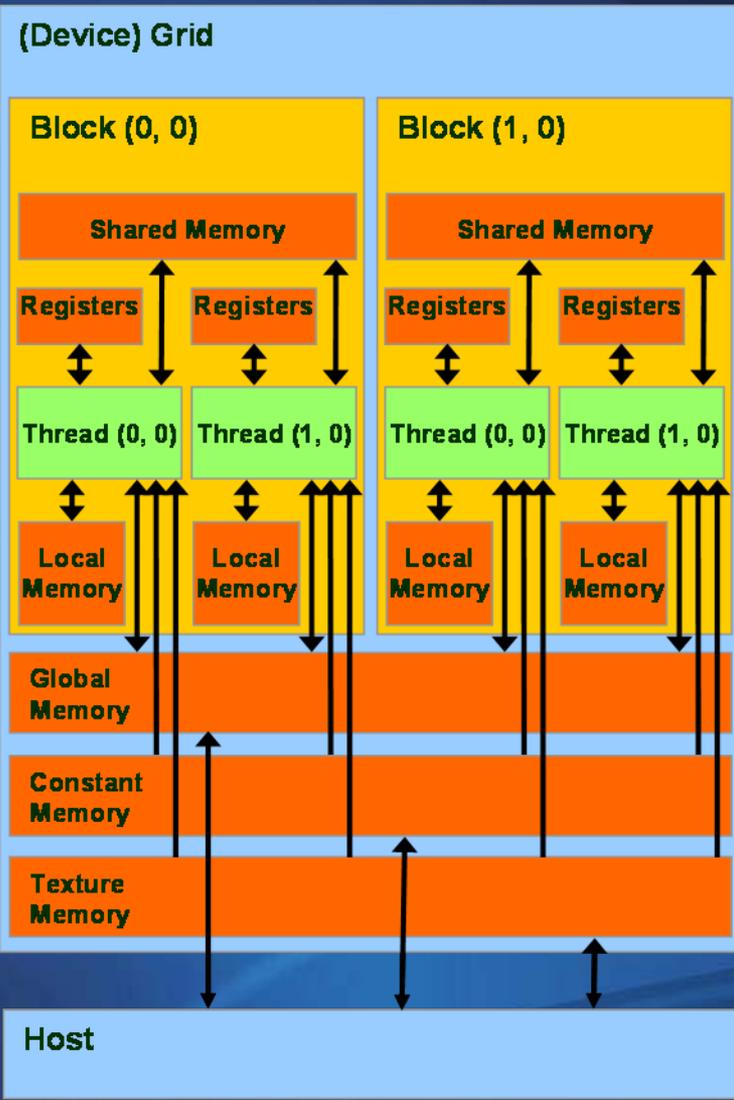


## NVIDIA GT200

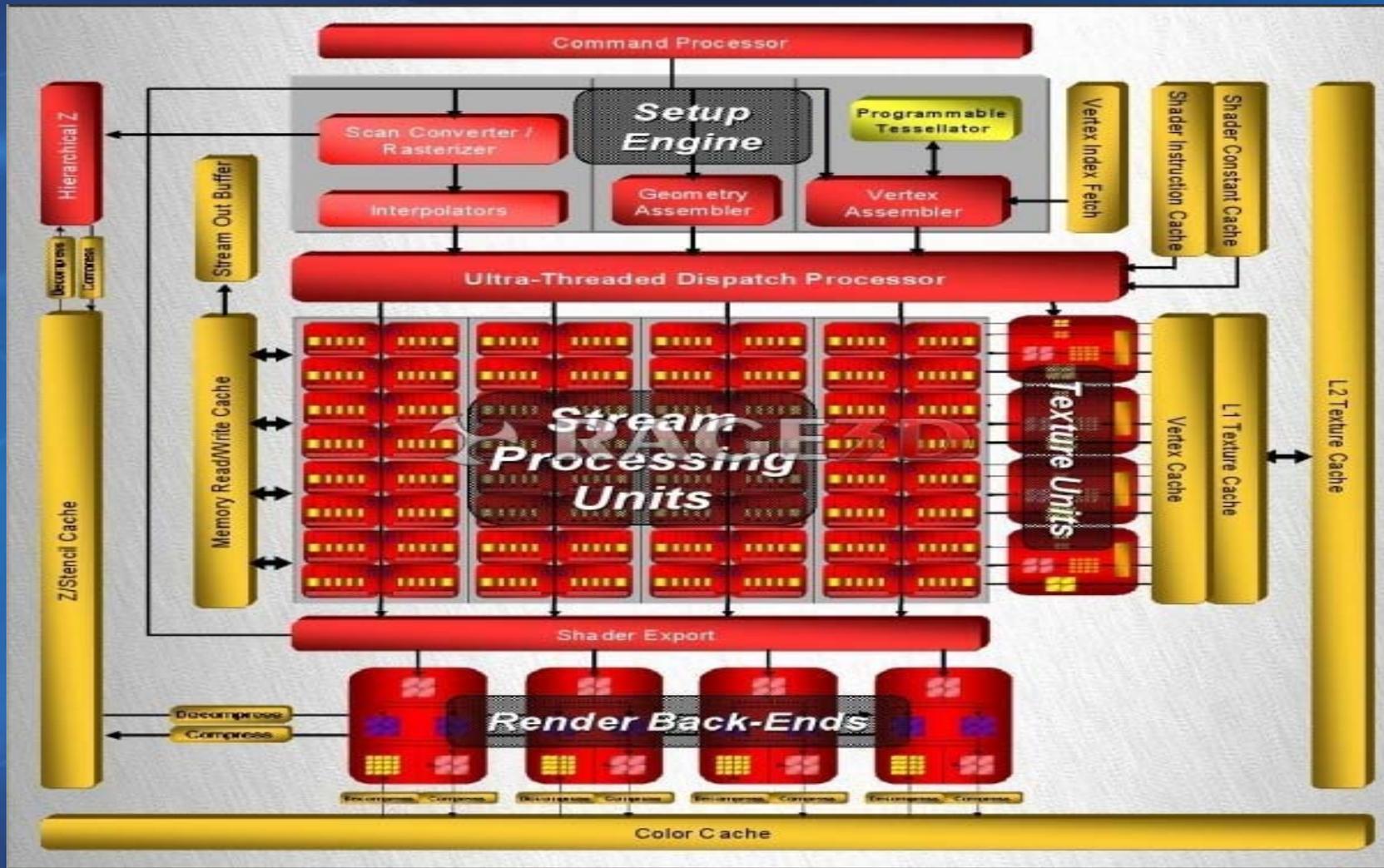
### Streaming Processor (SP)



# Nvidia 内存

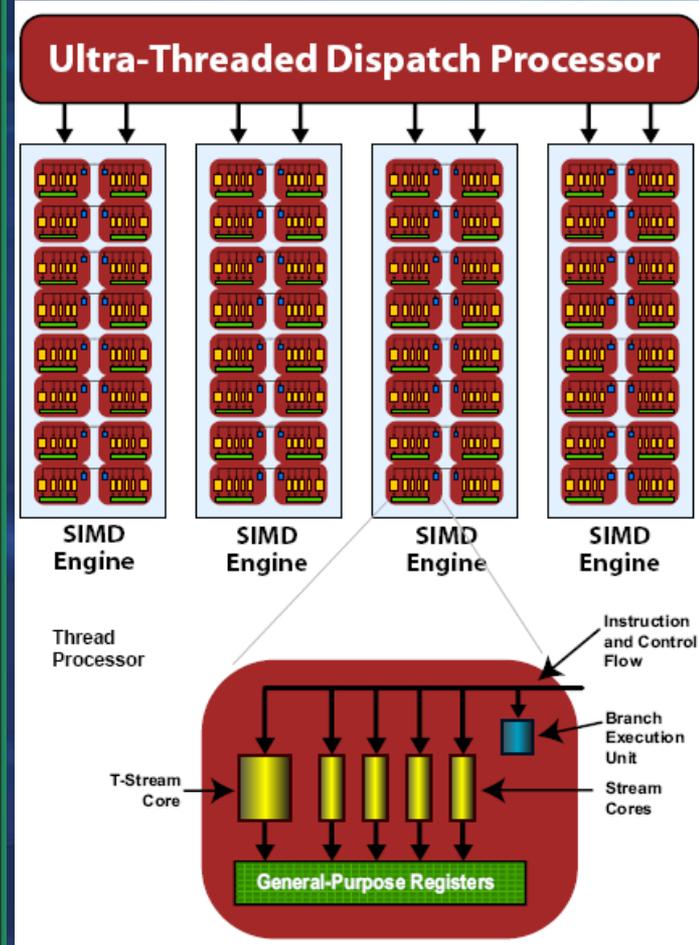
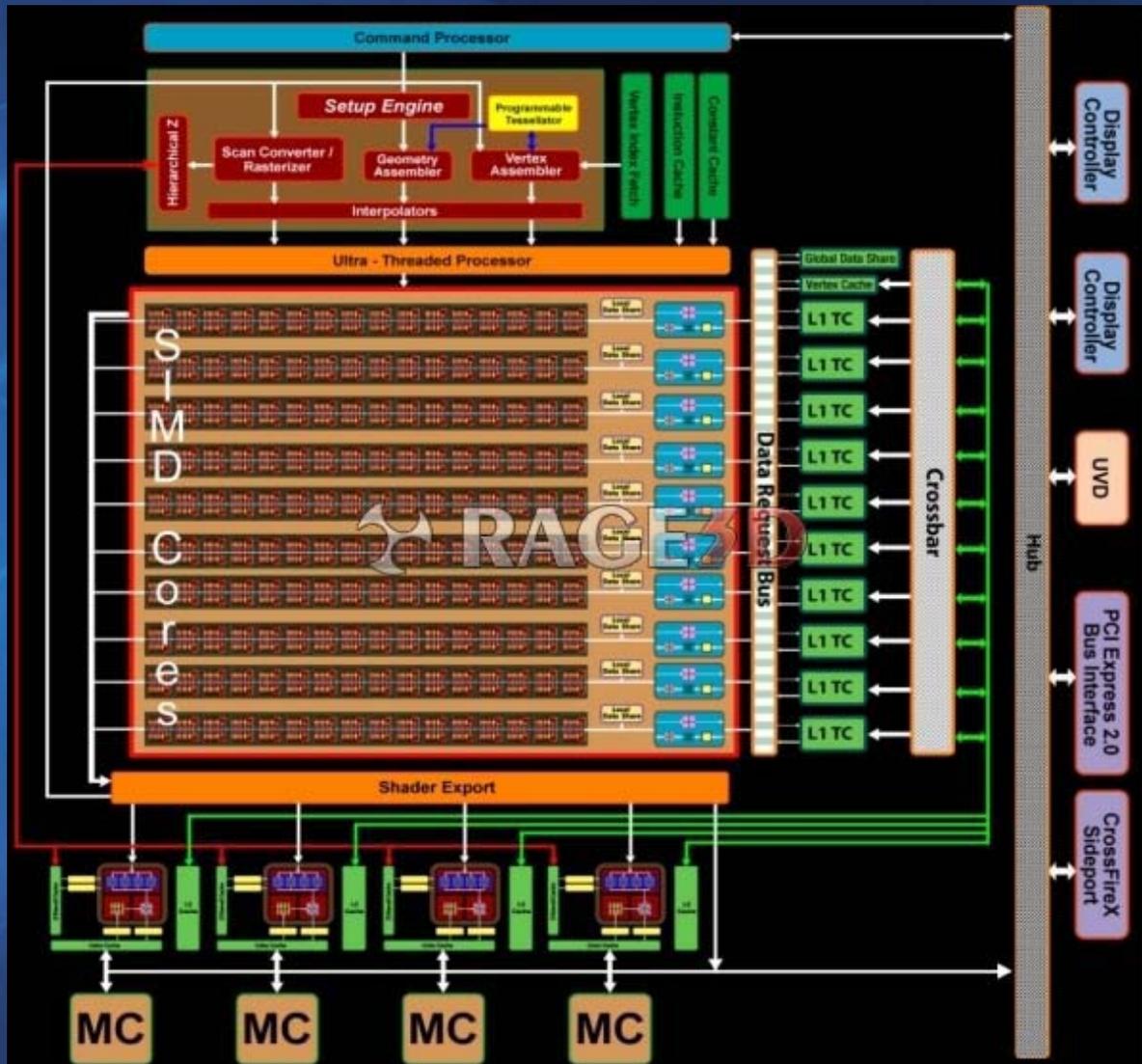


# AMD(ATI) R600, RV770硬件架构

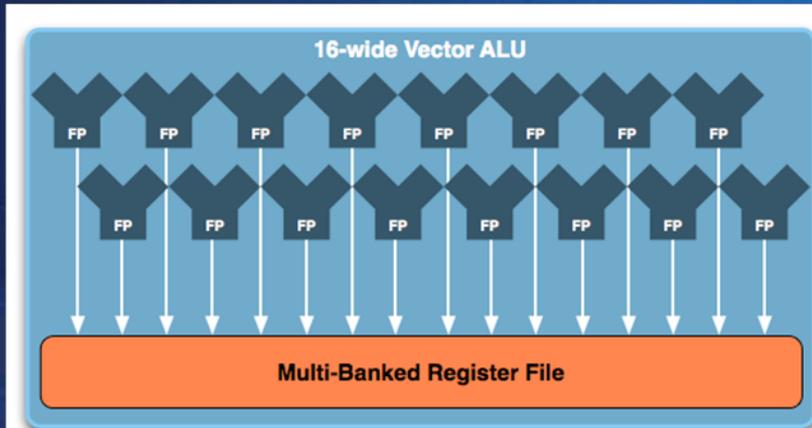


R600

# AMD(ATI) R600, RV770硬件架构

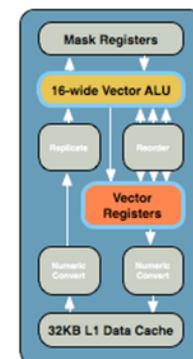
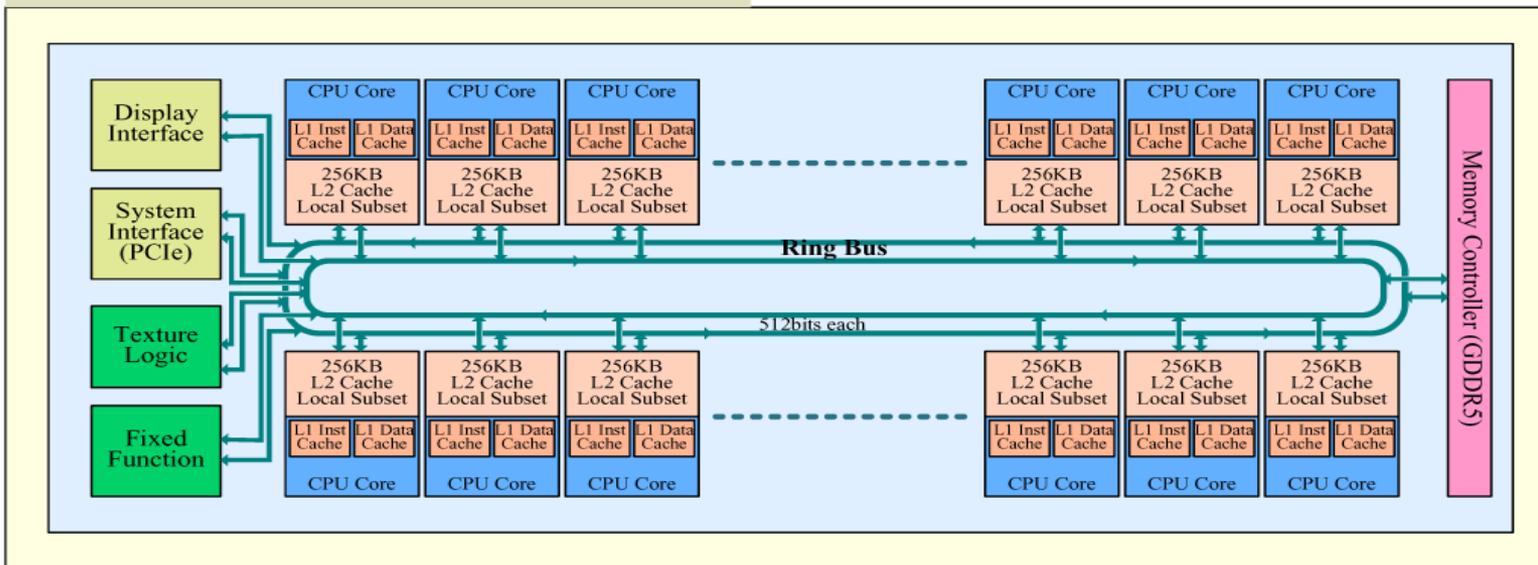


# Intel Larrabee架构



## Larrabee Block Diagram

Larrabee



# 提纲

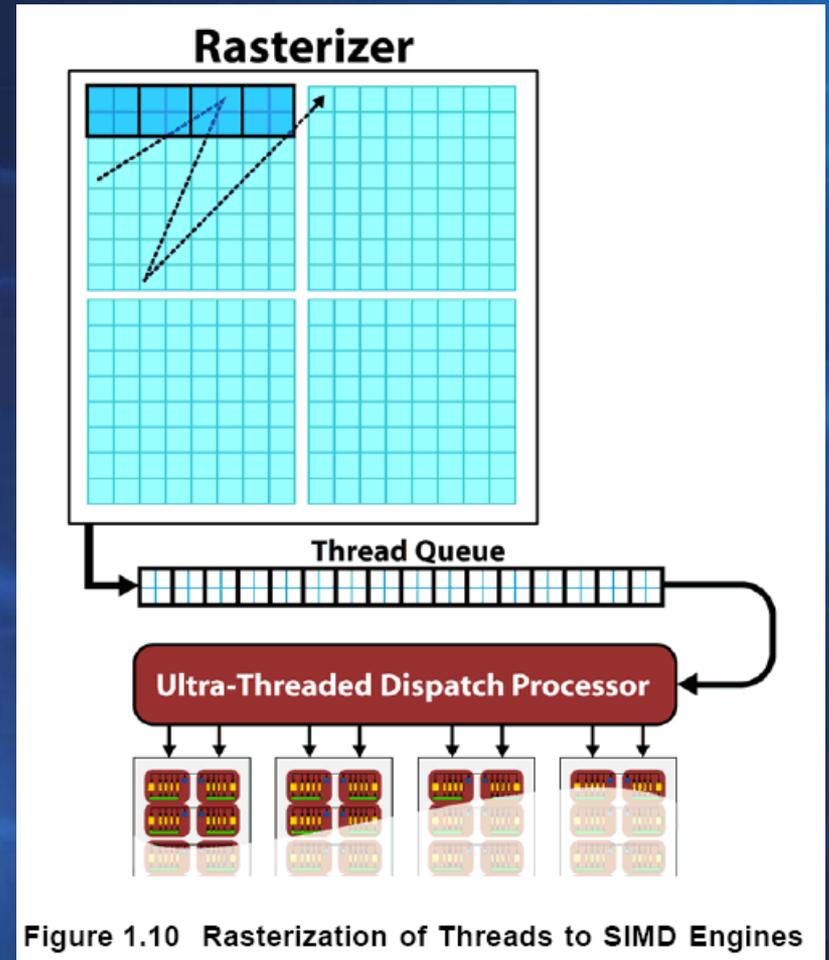
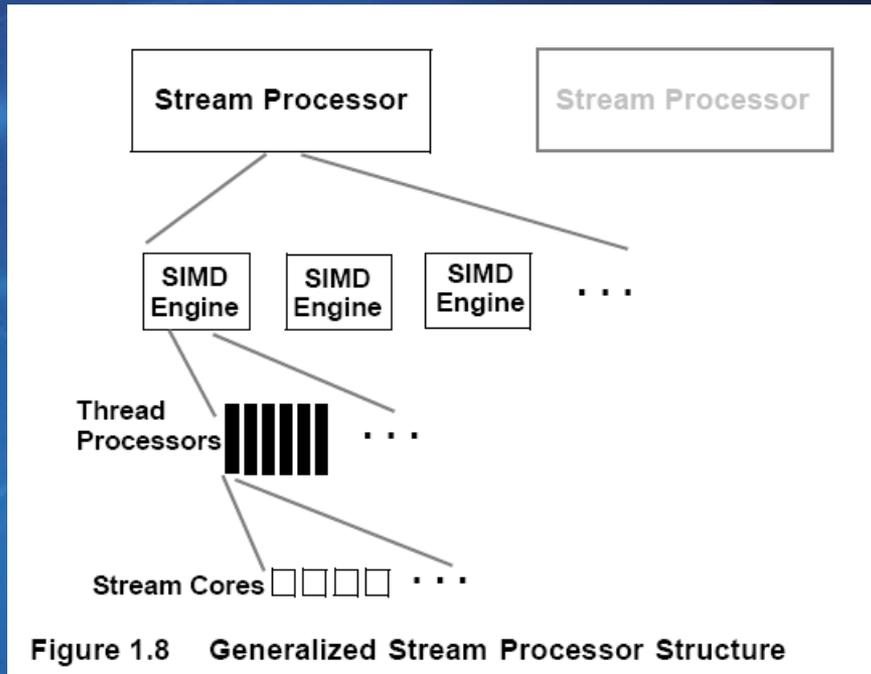
- 处理器核心技术趋势和现状
- 硬件架构简介
  - Nvidia: G80, GTX200
  - AMD(ATI): R600, RV770
  - Intel(Larrabee)
- 软件编程模型简介
  - Nvidia: CUDA
  - AMD(ATI): Firestream
  - Intel: Ct
- 应用与优化

# GPU技术路线的分析

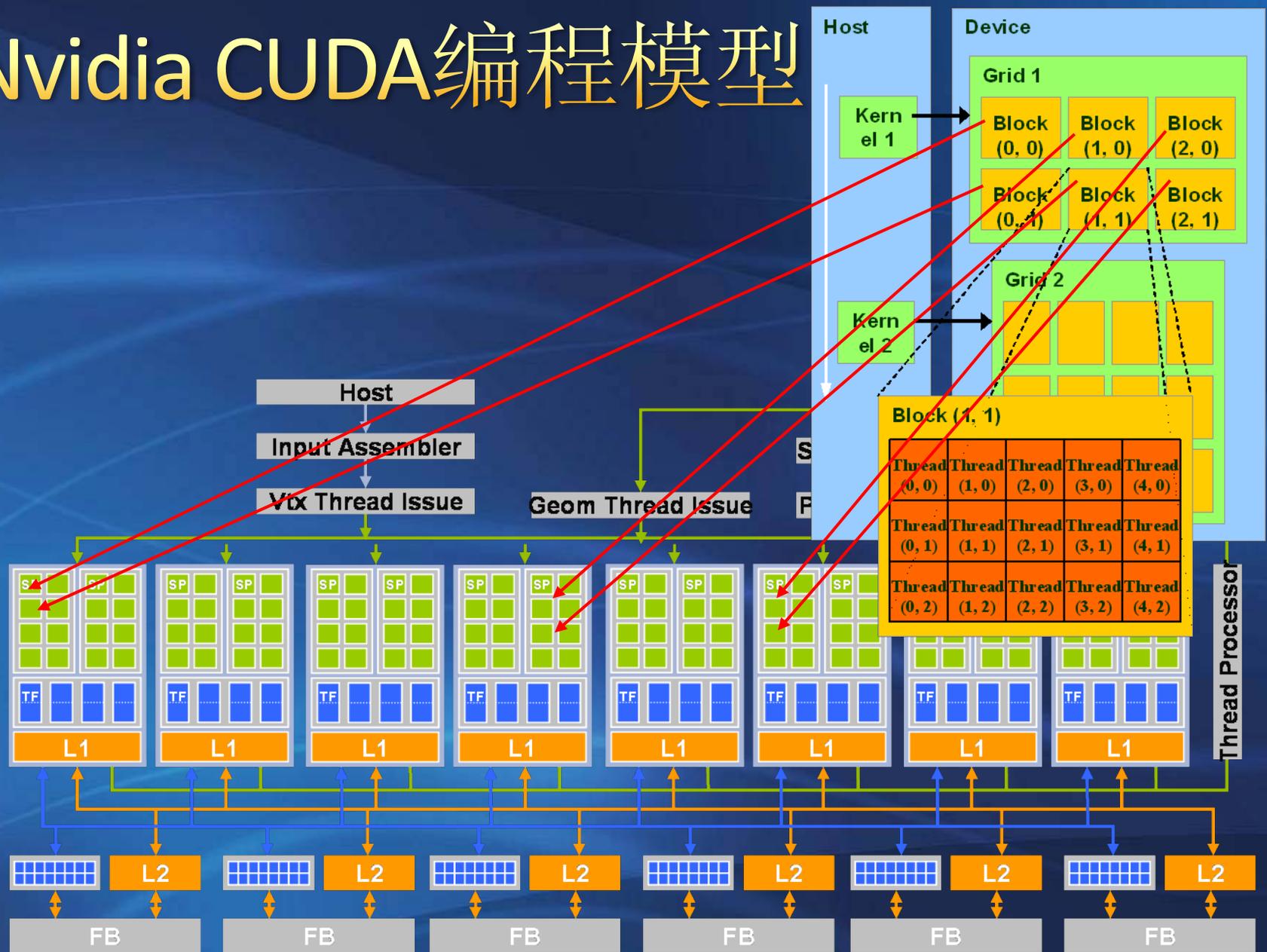
软件:

硬件	Nvidia(GT200)	AMD(RV770)	Intel(Larrabee)
X86	否	否	是
Cg	是	是	是
HLSL	是	是	是
ATI Stream	否	是	否
CUDA	是	否	否
OpenCL	是	是	是

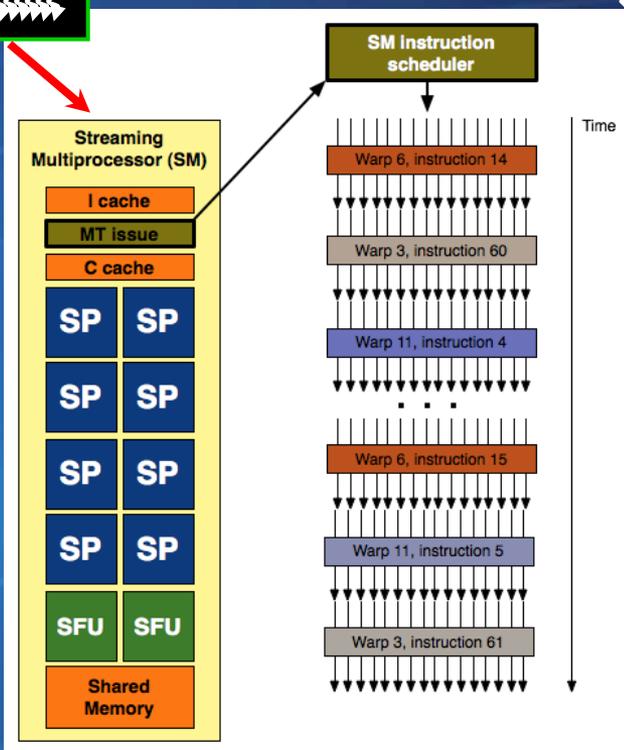
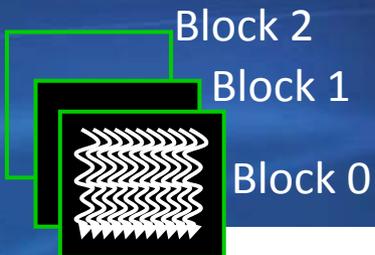
# AMD Stream computing



# Nvidia CUDA编程模型



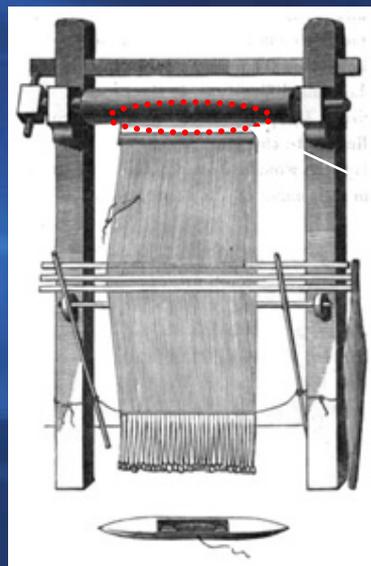
# CUDA线程执行模型



❖ SM内以(warp即32 threads)为单位并行执行

□ Warp(32 threads)内线程执行同一条指令

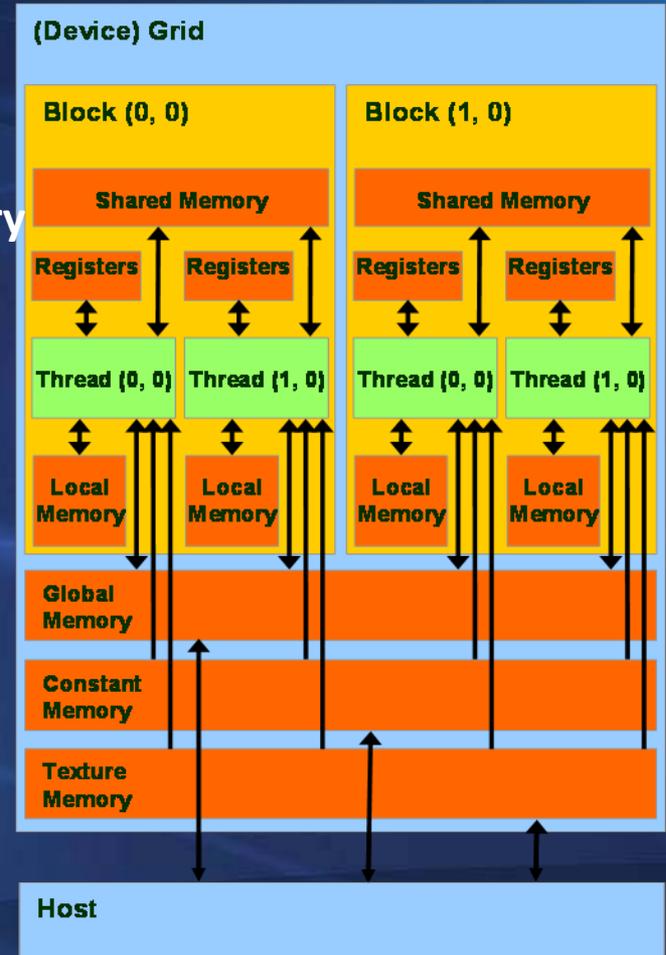
□ Half-warp(16 threads)是存储操作的基本单位



Warp

# CUDA 内存模型

- ❖ R/W per-thread **registers**
  - 1-cycle latency
- ❖ R/W per-thread **local memory**
  - Slow – register spilling to global memory
- ❖ R/W per-block **shared memory**
  - 1-cycle latency
  - But bank conflicts may drag down
- ❖ R/W per-grid **global memory**
  - ~500-cycle latency
  - But coalescing accessing could hide latency
- ❖ Read only per-grid **constant and texture memories**
  - ~500-cycle latency
  - But cached



# 提纲

- 处理器核心技术趋势和现状
- 硬件架构简介
  - Nvidia: G80, GTX200
  - AMD(ATI): R600, RV770
  - Intel(Larrabee)
- 软件编程模型简介
  - Nvidia: CUDA
  - AMD(ATI): Firestream
  - Intel: Ct
- 应用与优化

# GPU应用兴趣

- 09.3.2 发布浪潮倚天桌面超算，Intel Nehalem+Tesla GPU架构，国内客户表现出对GPU的极大兴趣
- 传统HPC领域：生物信息、遥感、计算化学、图像处理、天文研究、分子动力学、CAE、地震勘探、计算机体系研究、.....
- 新应用领域：3G、视频监控、网络安全、地图引擎、特种电影.....

# GPU应用开发合作模式

- 联合团队：客户应用团队+浪潮GPU开发团队
- 开发平台：浪潮倚天桌面超算、浪潮GPU集群
- 开发组织：算法并行化设计分析、移植
- 注意事项：热点定位、细节优化
- 案例：中科院基因组所+浪潮生物信息GPU开发

# GPU应用优化体会

- 项目并行性分析
  - 热点问题定位
  - 寻找数据并行性，减少数据之间的依赖性
  - 设计相应的并行算法
- 程序优化策略
  - 多的计算时间来隐藏内存访问延时
  - 合理分配线程和内存
  - 多用片上寄存器
  - 全局内存访问瓶颈的优化

QA?