



Performance Profiling

袁亚振



01 Introduction

02 Basic GPU Arch

03 GPU Profiling



- 帧率
 - FPS(Frames per second)
- 高帧率
 - 更流畅的画面
 - 更低的延迟
- 帧率要求
 - 平台/游戏类型
 - 主机: 60FPS
 - 手机: 30FPS?
 - VR : 90FPS+
 - 平稳



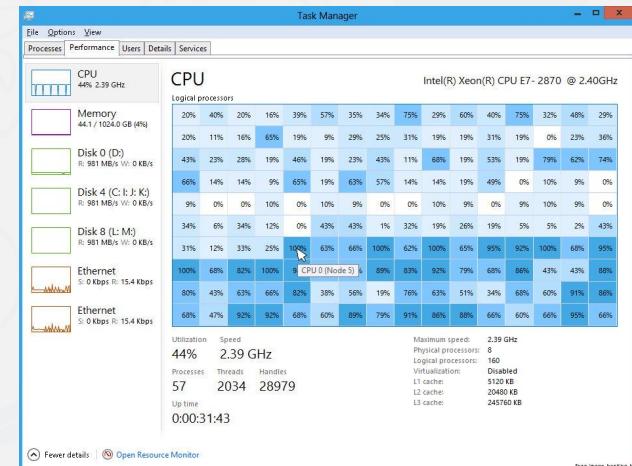
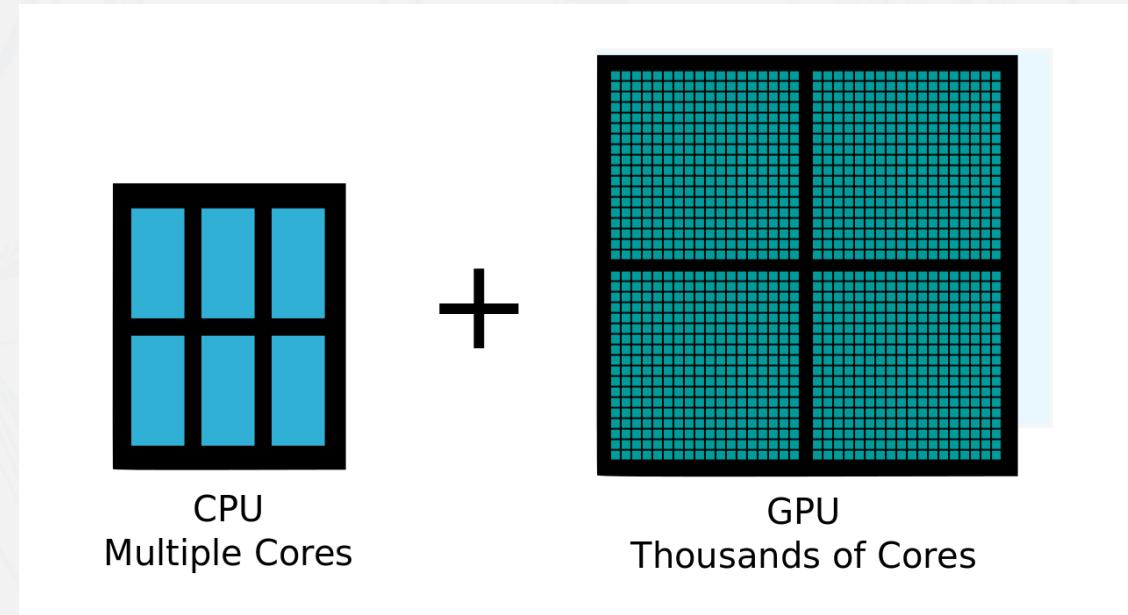


- CPU 与 GPU 架构的不同

CPU	GPU
<ul style="list-style-type: none">* Low compute density* Complex control logic* Large caches (L1\$/L2\$, etc.)* Optimized for serial operations<ul style="list-style-type: none">• Fewer execution units (ALUs)• Higher clock speeds* Shallow pipelines (<30 stages)* Low Latency Tolerance* Newer CPUs have more parallelism	<ul style="list-style-type: none">* High compute density* High Computations per Memory Access* Built for parallel operations<ul style="list-style-type: none">• Many parallel execution units (ALUs)• Graphics is the best known case of parallelism* Deep pipelines (hundreds of stages)* High Throughput* High Latency Tolerance* Newer GPUs:<ul style="list-style-type: none">• Better flow control logic (becoming more CPU-like)• Scatter/Gather Memory Access• Don't have one-way pipelines anymore

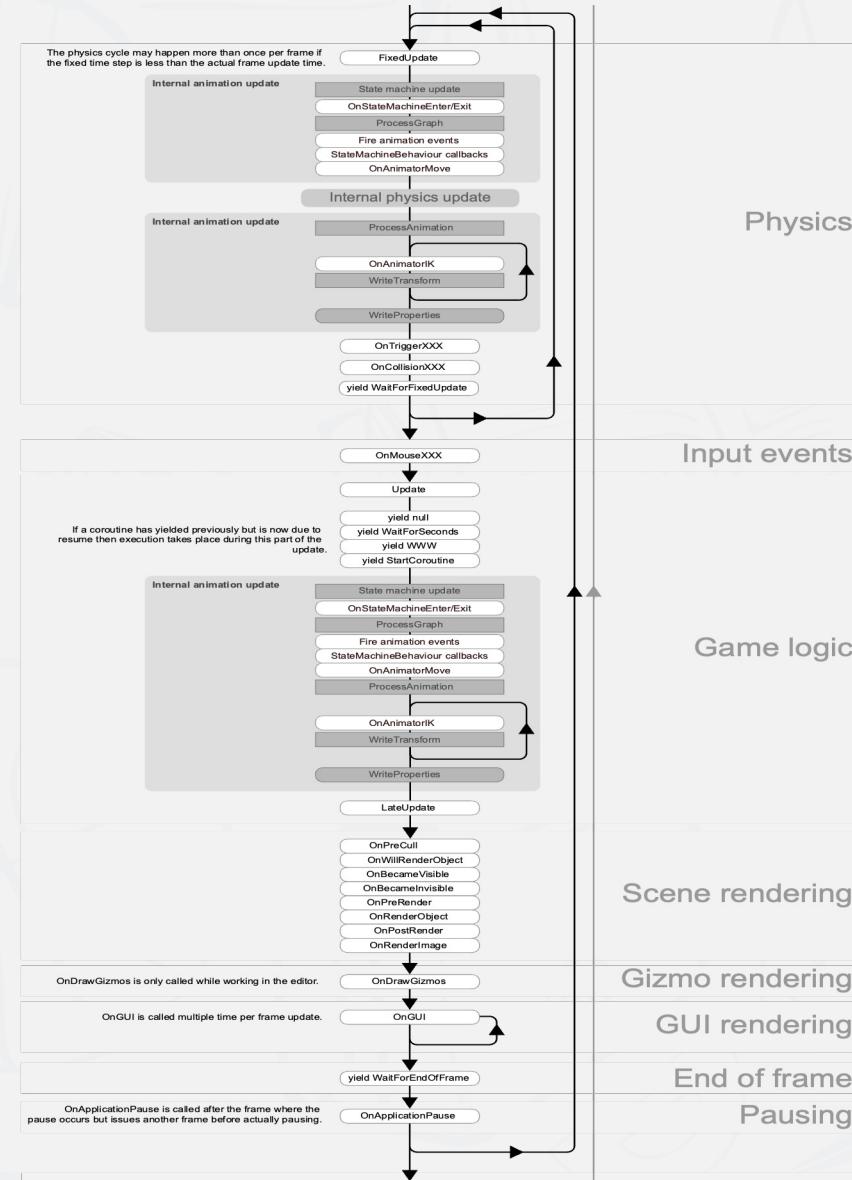


- CPU 和 GPU 需要很好的协作
 - “Batch, Batch, Batch” [GDC2003]
 - CPU的核数也在增长



Game Loop

Introduction





Gamethread

Frame N Frame N+1 Frame N+2

RenderThread (VkCmds)

Frame N Frame N+1 Frame N+2

GPU (double buffer)

Frame N Frame N+1 Frame N+2

Display

Frame N Frame N+1 Frame N+2

Gamethread vs Renderthread

Introduction



Gamethread

Frame N Frame N+1 Frame N+2

RenderThread (VkCmds)

Frame N Frame N+1 Frame N+2

GPU (double buffer)

Frame N Frame N+1 Frame N+2

Display

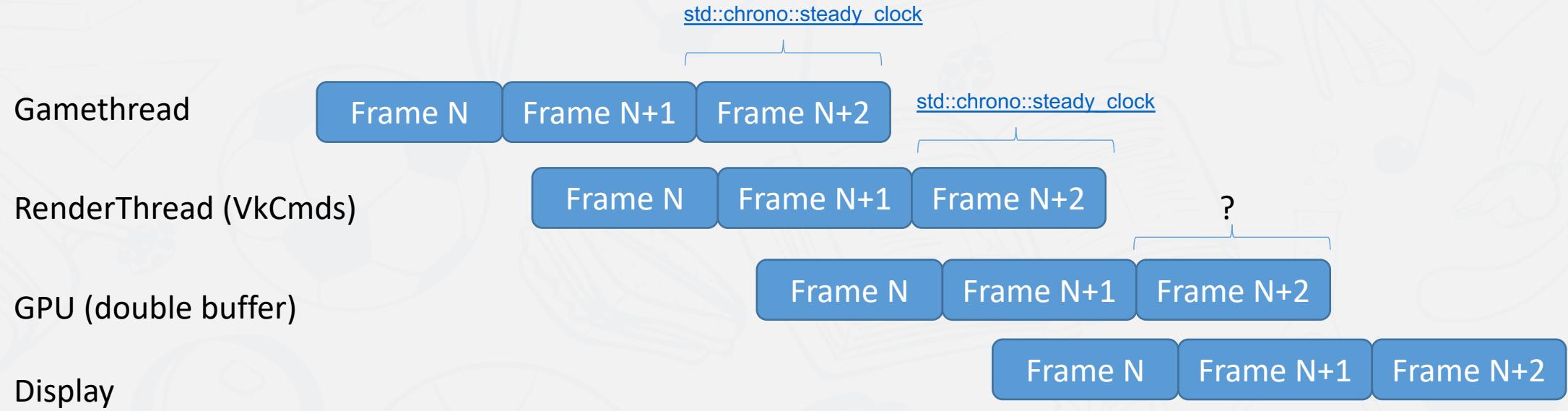
Frame N Frame N+1 Frame N+2

最慢的一环决定了最终的FPS



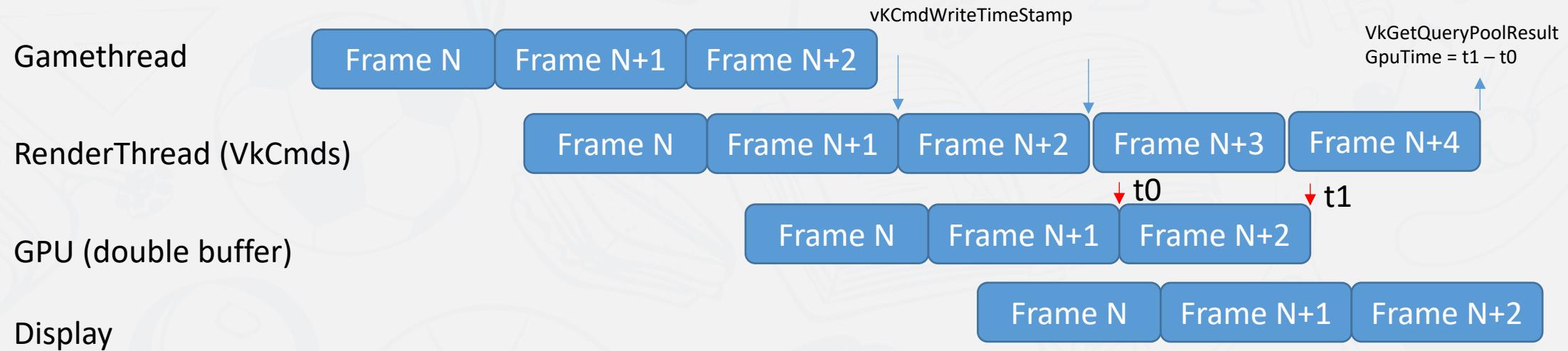


- Measure Time





- Measure Time



```
// Provided by VK_VERSION_1_0
void vkCmdWriteTimestamp(
    VkCommandBuffer
    VkPipelineStageFlagBits
    VkQueryPool
    uint32_t
        commandBuffer,
        pipelineStage,
        queryPool,
        query);
```

```
// Provided by VK_VERSION_1_0
VkResult vkGetQueryPoolResults(
    VkDevice
    VkQueryPool
    uint32_t
    uint32_t
    size_t
    void*
    VkDeviceSize
    VkQueryResultFlags
        device,
        queryPool,
        firstQuery,
        queryCount,
        dataSize,
        pData,
        stride,
        flags);
```



GameThread

- 脚本
- 物理
- 网络
- ...

RenderThread

- 动态物体过多
- 灯的数量过多
- DrawCall 过多
- ...

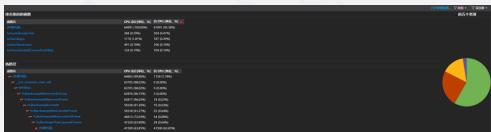
GPU

- 三角形数量
- 贴图数量
- 分辨率
- Raytracing 过于复杂
- 显存
- ...



GameThread

- 脚本
- 物理
- 网络
- ...



Visual studio Profiler

RenderThread

- 动态物体过多
- 灯的数量过多
- DrawCall 过多
- ...

GPU

- 三角形数量
- 贴图数量
- 分辨率
- Raytracing 过多
- 显存
- ...

RenderDoc
Nsight

Pix

Intel® VTune™ Profiler

Find and Fix Performance Bottlenecks Quickly and Realize All the Value of Your Hardware



Basic GPU Arch

GPU Architecture

Basic GPU Arch

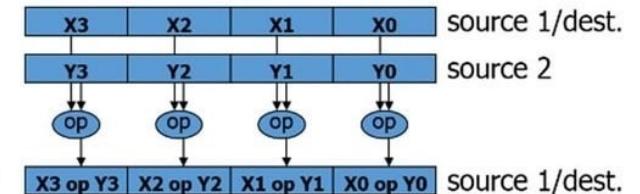


CPU

SIMD

1 instruction – multiple data

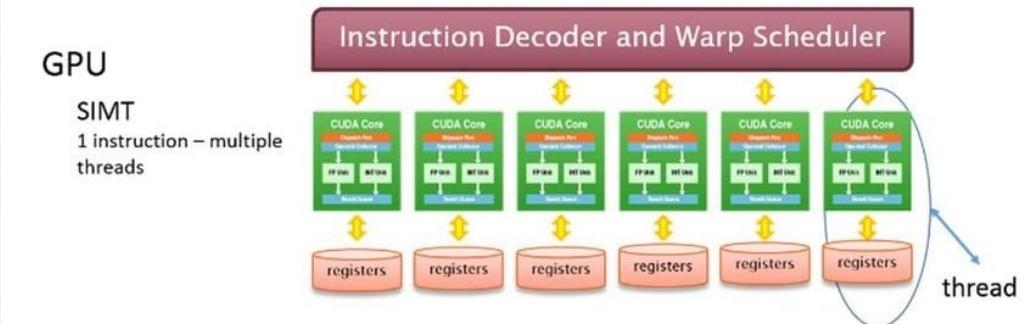
SSE2/3/4 – Neon – Altivec
AVX – AVX2...



GPU

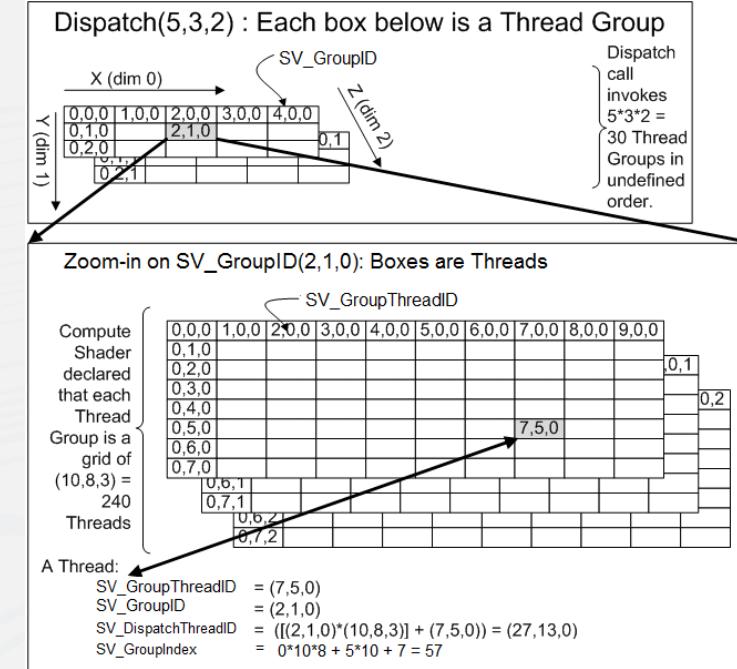
SIMT

1 instruction – multiple threads



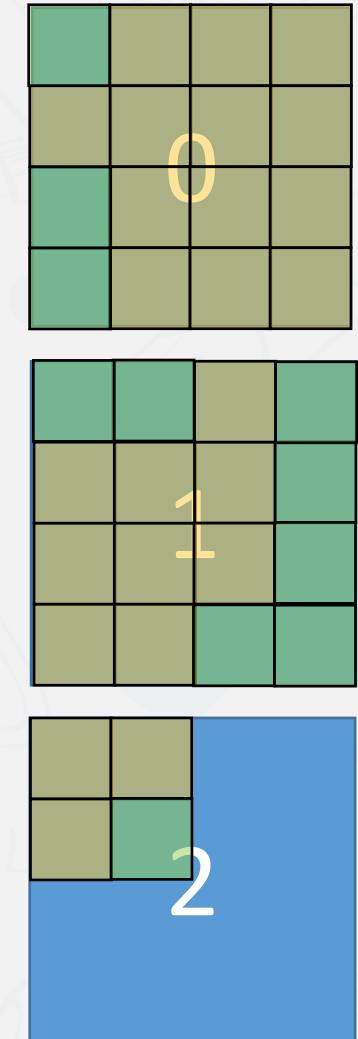
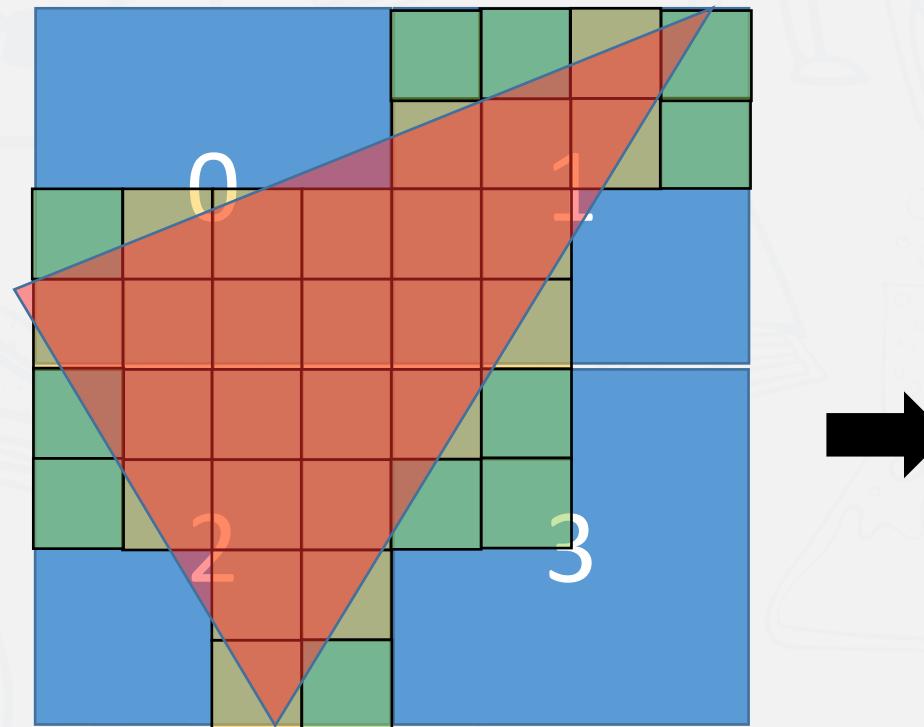


- Warp & Wavefront
 - 一次指令发射多个线程一起执行
 - Warp : 32 线程 , NVIDIA
 - Wavefront : 64线程 AMD
- SIMT 执行
 - 编程者只用关心自己的线程
- Compute Shader
 - 线程组中warp数量 : $(\text{thread} + 32 - 1) / 32$
 - Numthread(10,8,3) : 8 warp
 - Inactive Thread/lanes: $8 * 32 - 240 = 16$





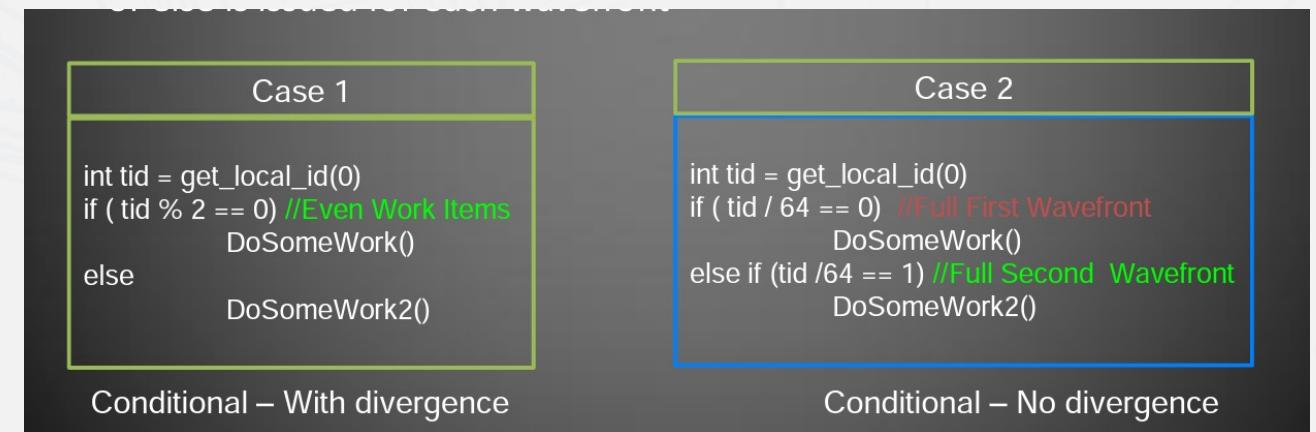
- Warp & Wavefront
- SIMT 执行
- Compute Shader
- Vertex Shader
 - Batch Processing
 - Post-Transform Vertex Cache
- Pixel Shader
 - Quad 为单位
 - Tile Walk & PS Wave pack
 - Helper lane/inactive thread



[Drobot. et al 2020]

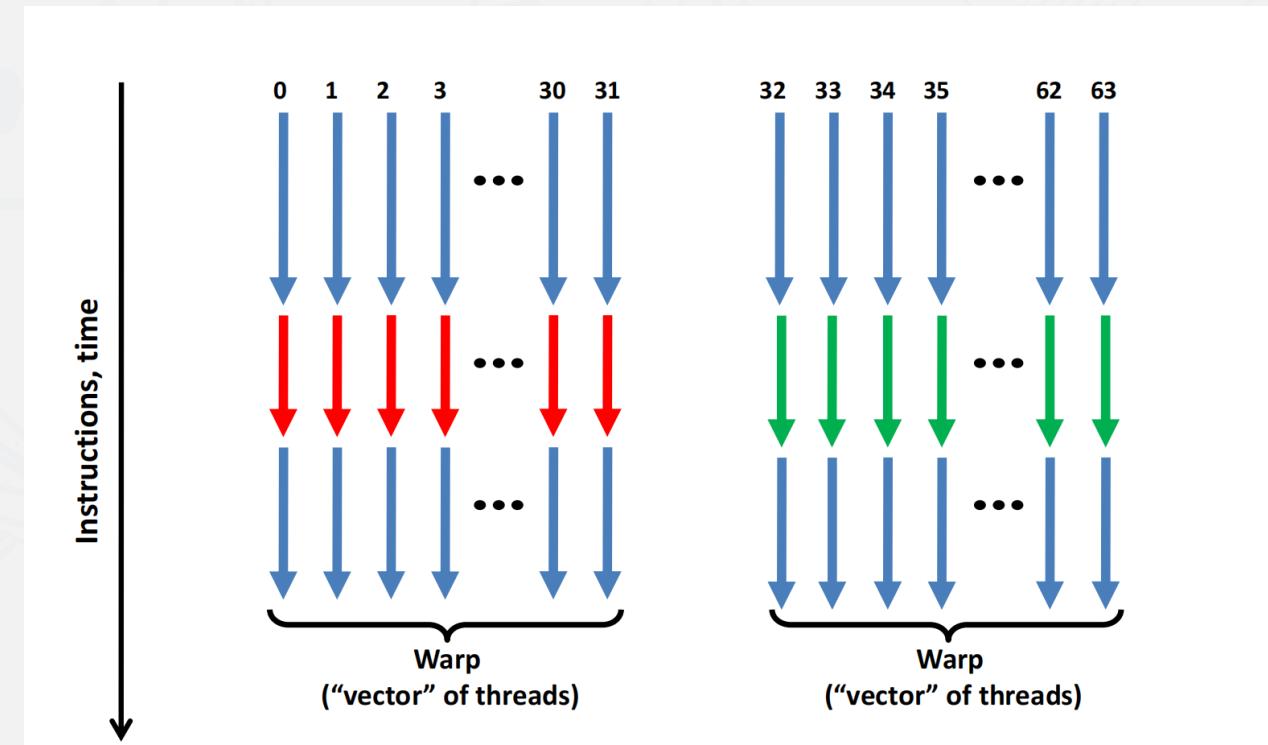


- 以Warp为单位执行指令
 - Case 1: Warp 内走的是不同分支
 - Case 2: Warp 内走的是一样的分支





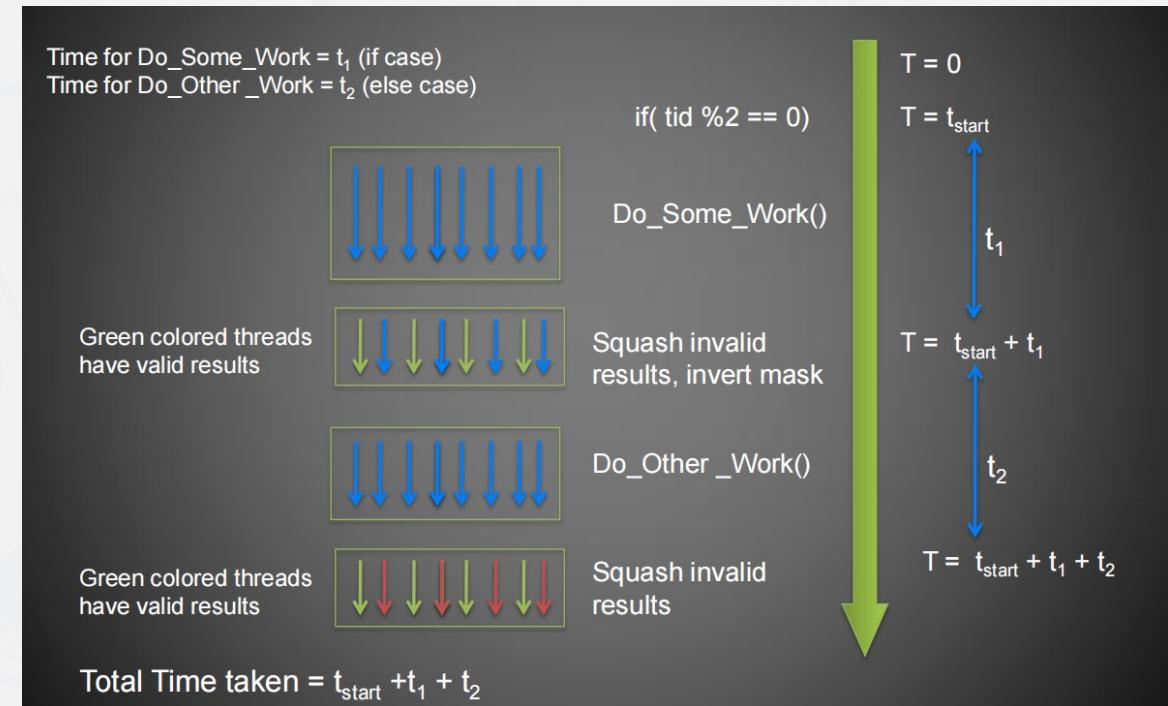
- Case 1
 - 没有性能负担
 - 比如uniform



Case 1



- Case 2
 - 两个分支都会走
 - 算术指令都会执行
 - 访存: 读和写都不会发生

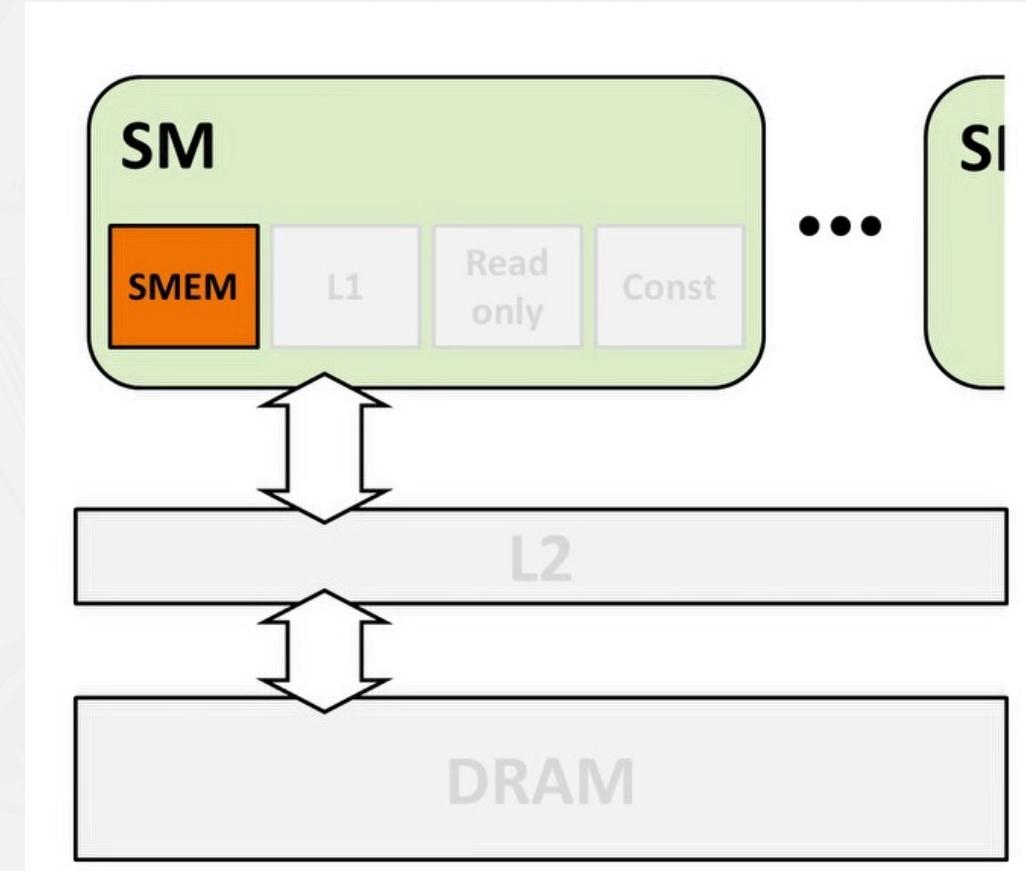


- 4个Warp Scheduler , 双发射
- Instruction Latency
- Warp Switching
 - 没有Overhead
- 寄存器/SMEM资源是预先分配的
 - Threadblock
 - Occupancy



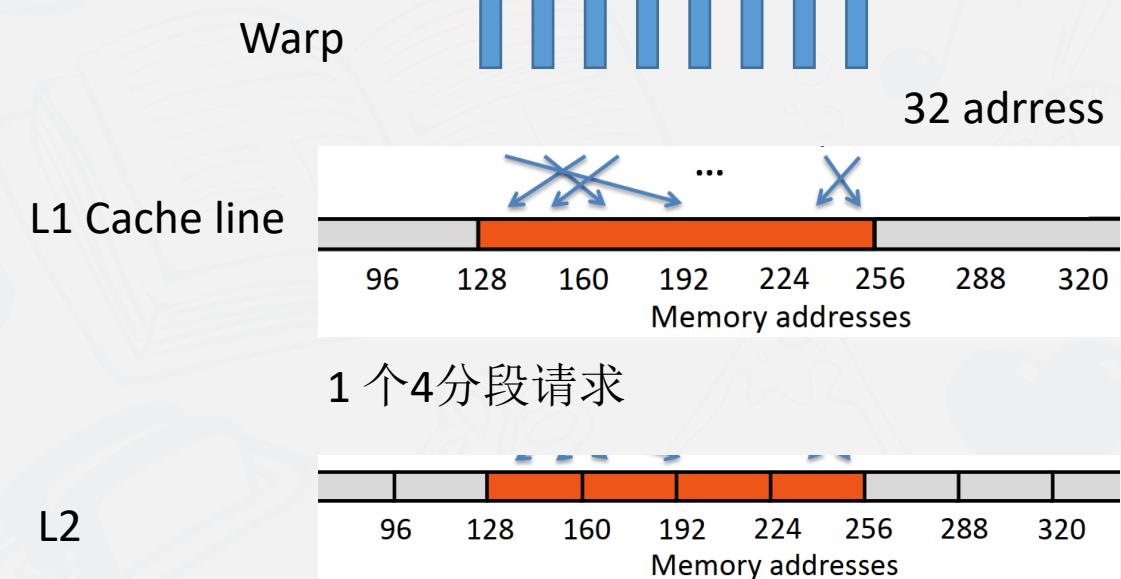


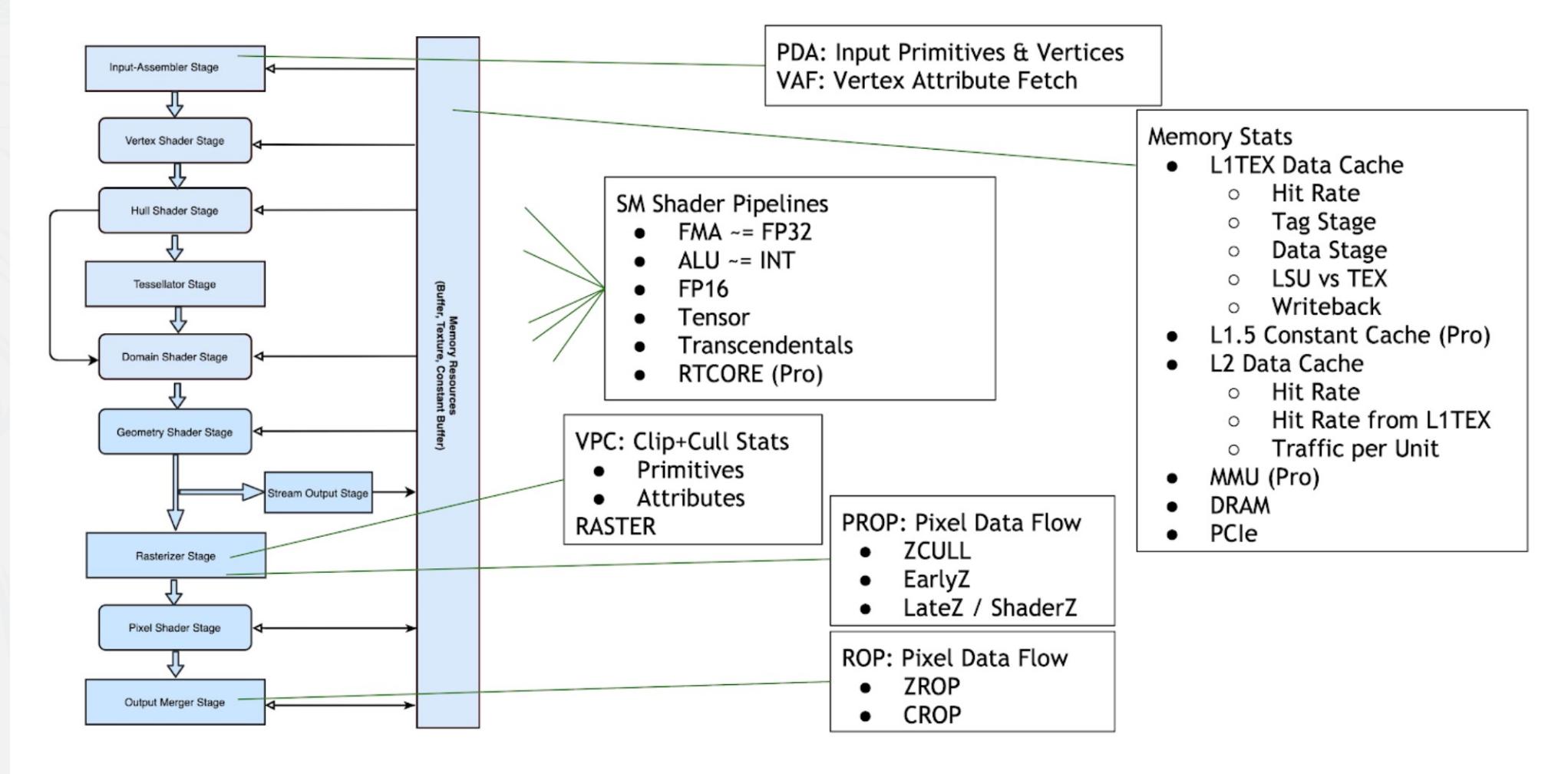
- DRAM
 - Read /write to Video RAM
- L2
- L1
- Read Only
 - Textures
- Const
 - uniforms
- Shared Memory





- Cache Line :128 bytes
- Segments: 32bytes
- Warp 提供32个地址
 - 硬件整合成请求
 - Texture是4个thread为单位



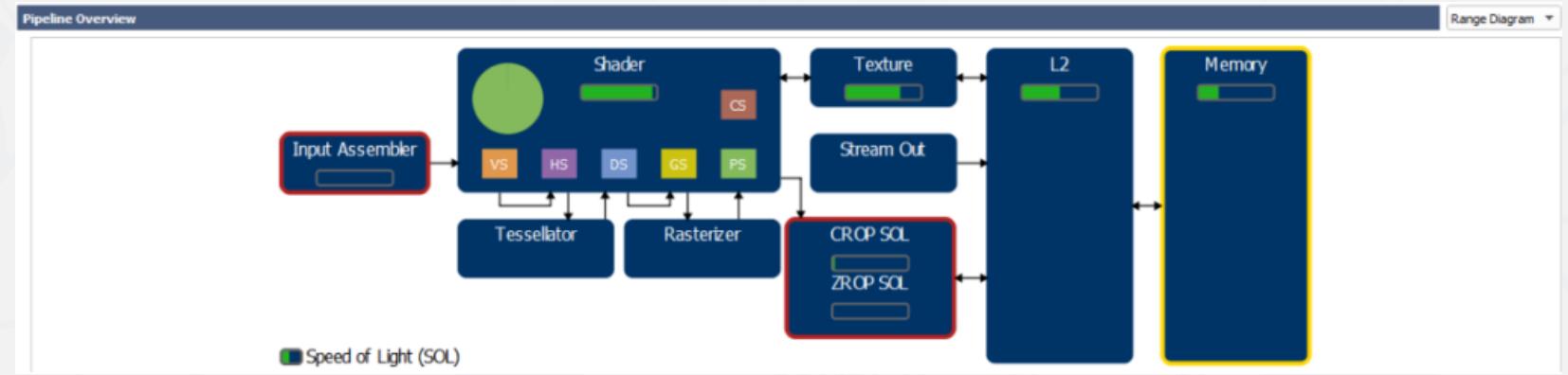




Nsight GPU Profiling



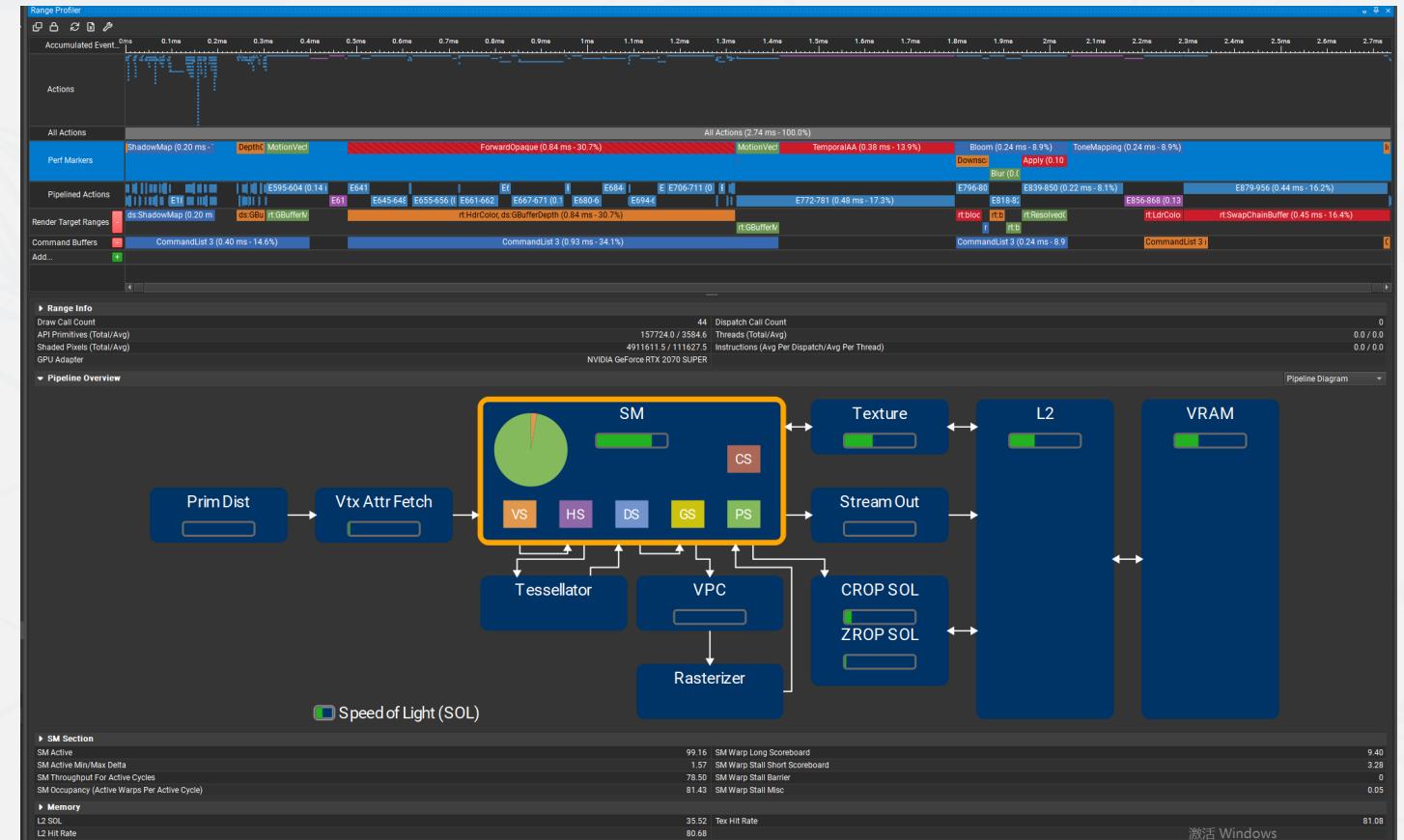
- Nvidia GPU Metrics
- Throughput based
 - SOL



"The Peak-Performance-Percentage Analysis Method for Optimizing Any GPU Workload"

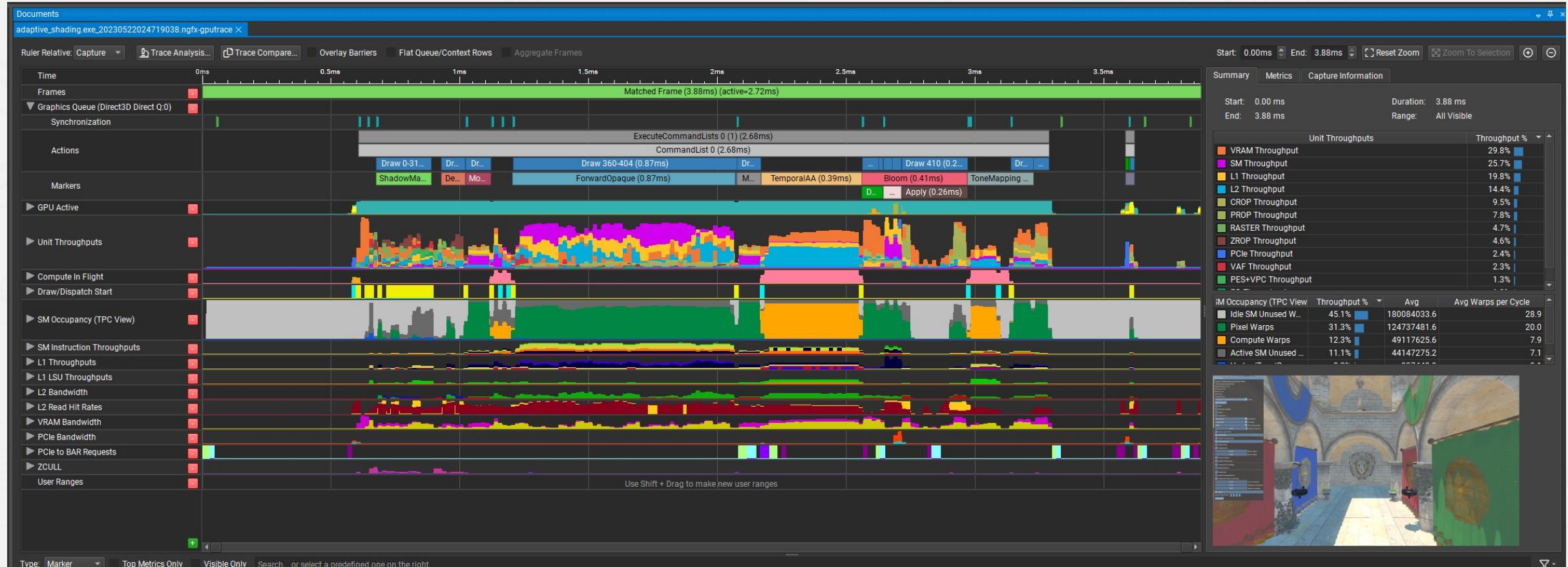


- 基于指令重放
- 可以分析每一个GPU Event的具体情况



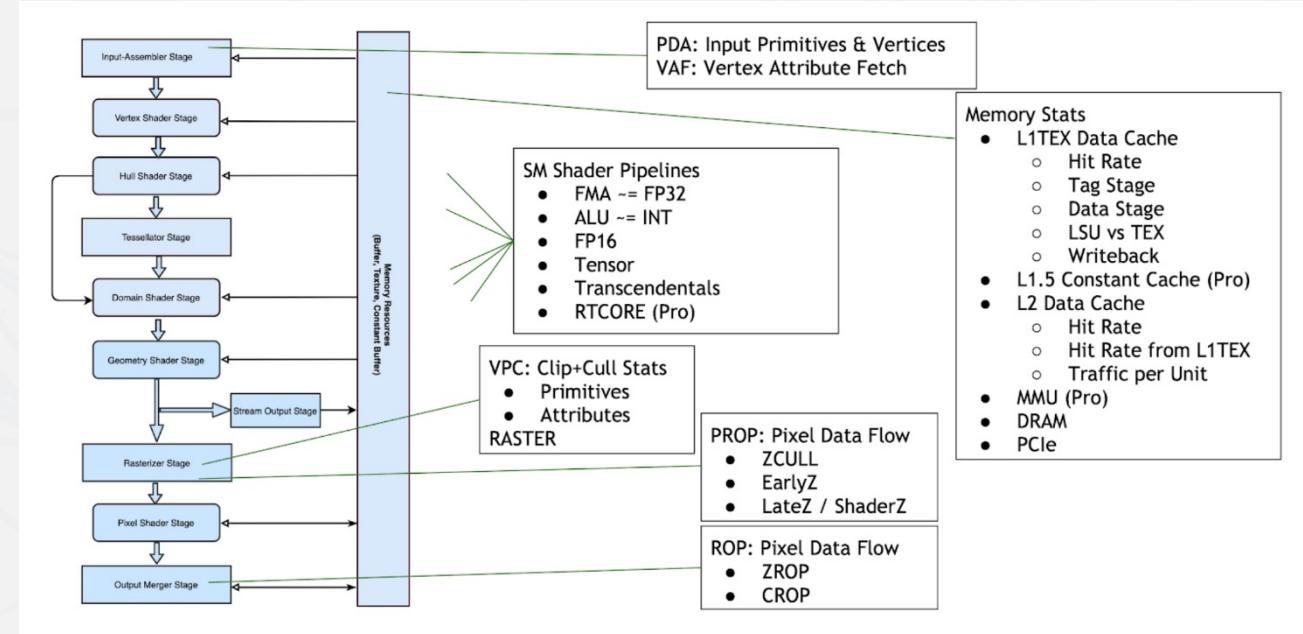


- 即时的结果，实时捕获，时间值上更准确。





- VAF
- PD
- PES+VPC
- RASTER
- SM
- L1TEX
- L2
- VRAM
- ZROP
- CROP



SM Occupancy

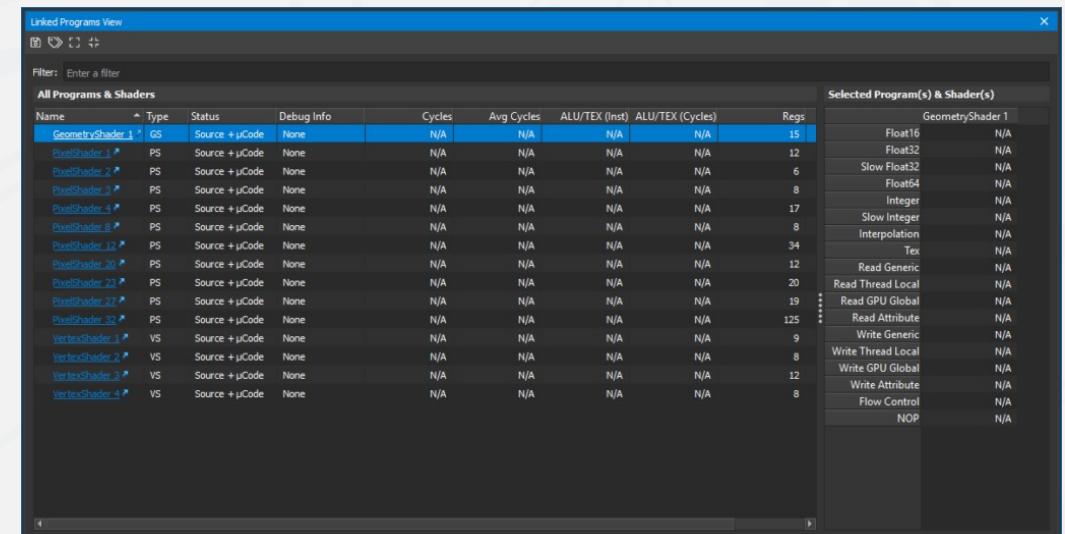
Nsight GPU Profiling



- 一个SM上最多活跃多少个Warp
 - Hide latency
- Register Number
 - GPU Profiler->Linked Program



SM Occupancy



Impact of Varying Register Count Per Thread

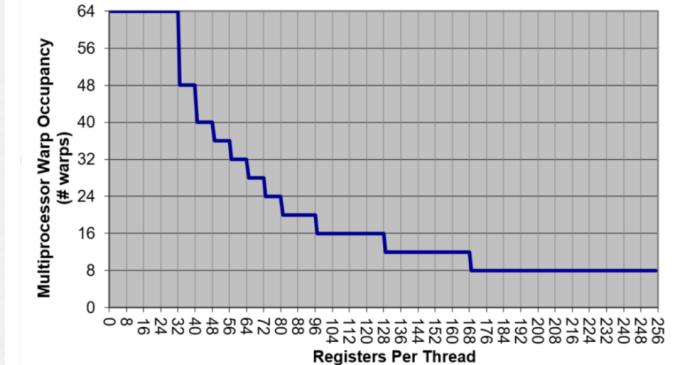
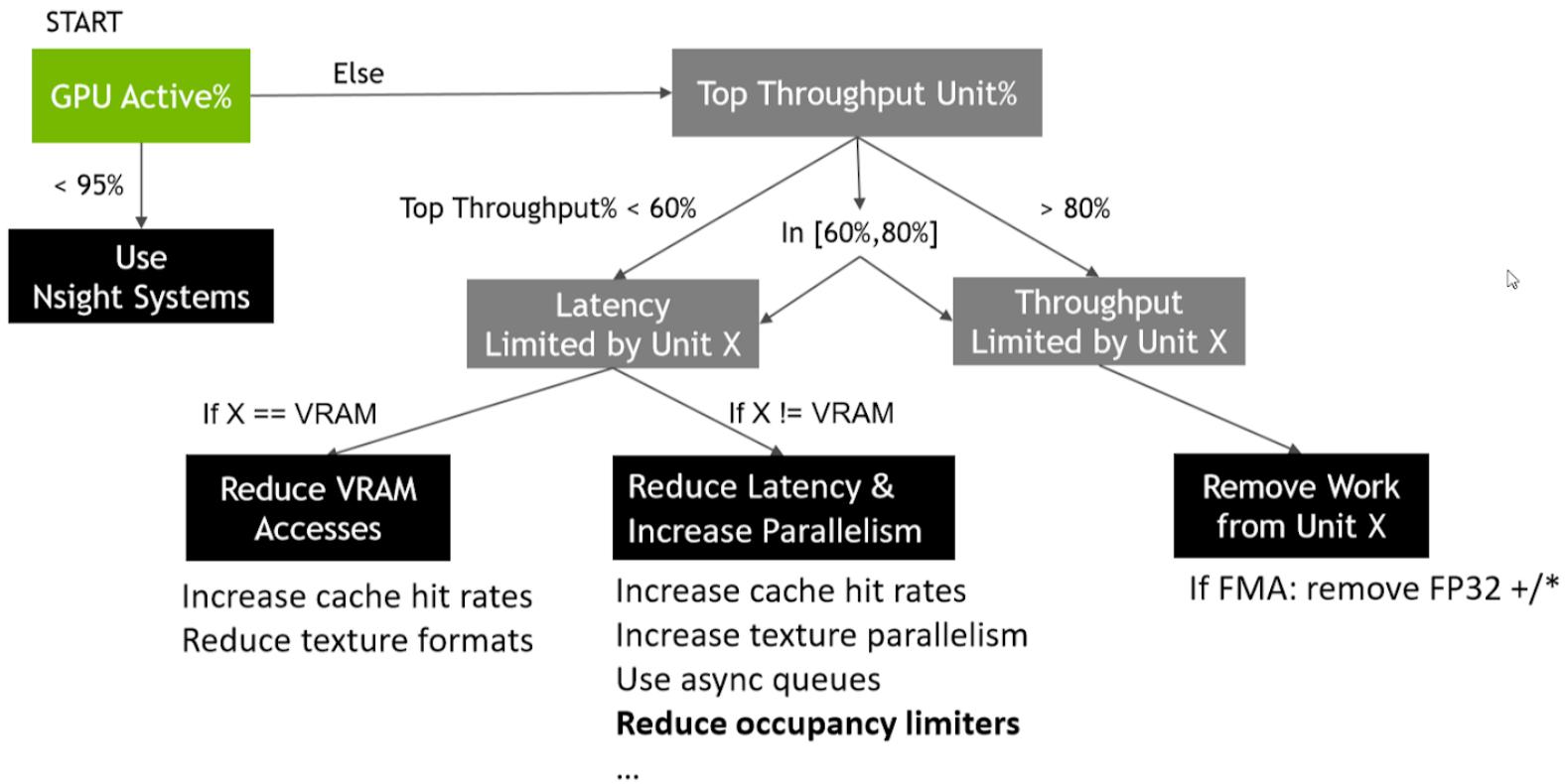


Figure 13. Graph from the CUDA Occupancy Calculator for "Compute Capability" 6.1.



THE P3 (PEAK-PERF%) METHOD



不同的平台瓶颈可能不一样，以实测为基准！



- VAF
- PD
- PES+VPC
- RASTER
- SM
- L1TEX
- L2
- VRAM
- ZROP
- CROP

Geometric LOD

Multi Rate Shading

Shader LOD

Texture LOD

Reference



- <https://developer.nvidia.com/blog/the-peak-performance-analysis-method-for-optimizing-any-gpu-workload/>
- <https://on-demand.gputechconf.com/gtc/2013/presentations/S3466-Programming-Guidelines-GPU-Architecture.pdf>
- <https://docs.nvidia.com/nsight-graphics/UserGuide/index.html>