# ECE/CS 472/572:
# Special Topics Part III:
# GPU Architecture Introduction

Prof. Lizhong Chen

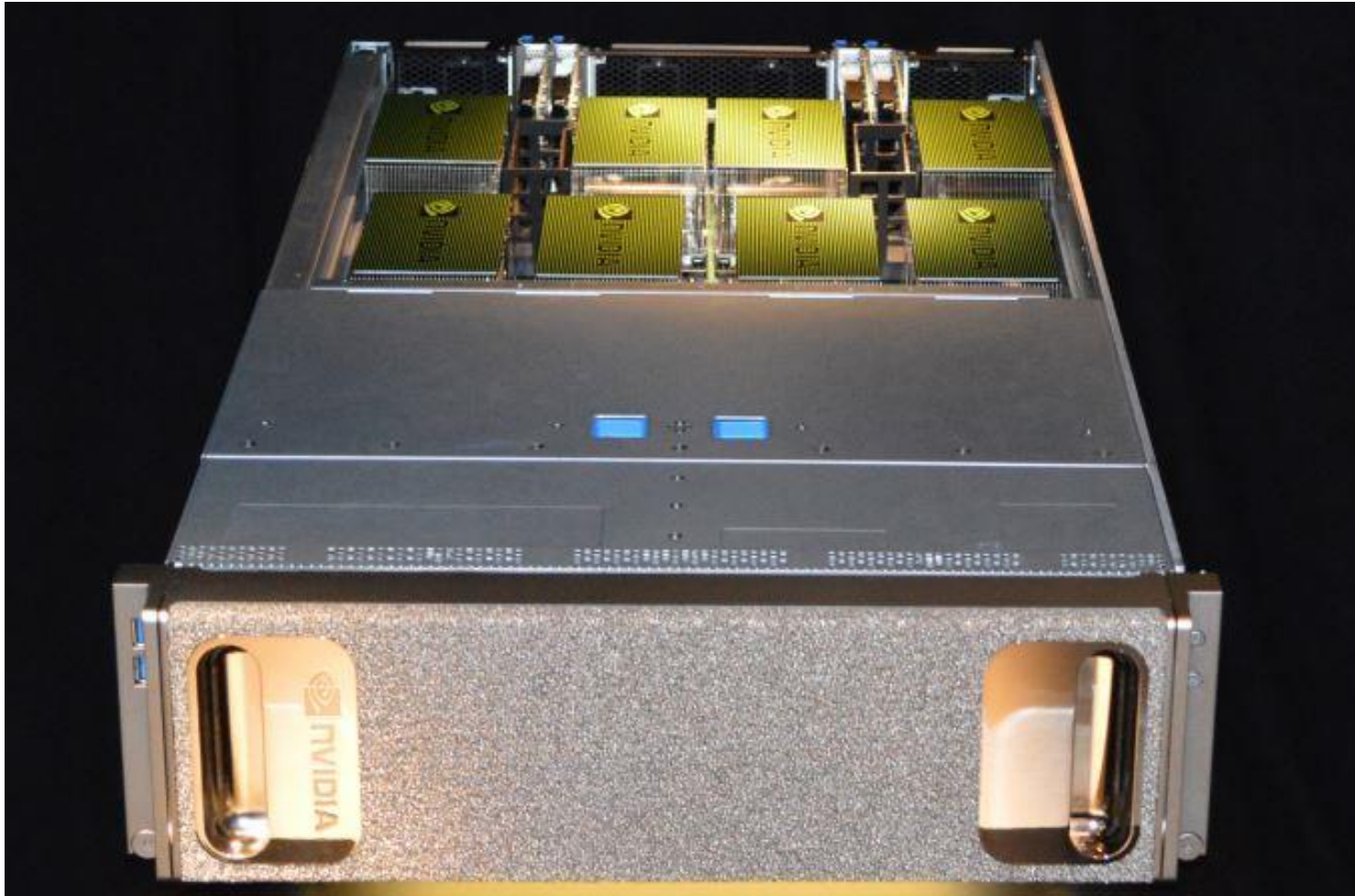Spring 2019

# GPU: System Organization

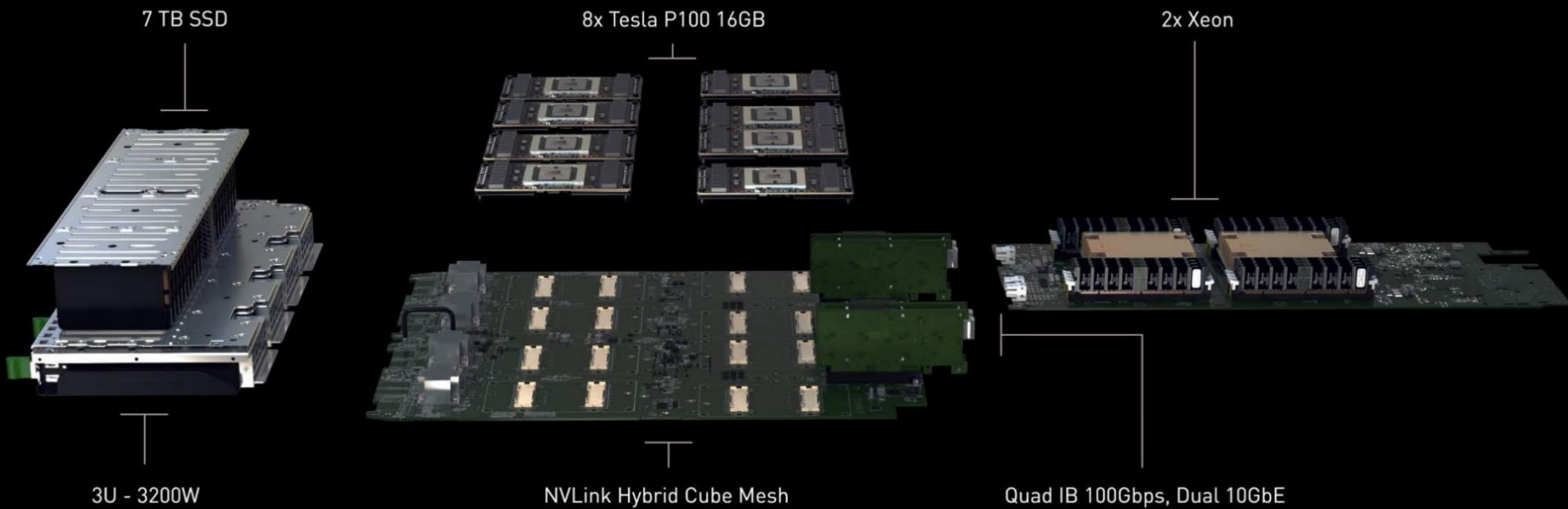- Nvidia DGX-1 High-performance computing (HPC) sever

# GPU: System Organization
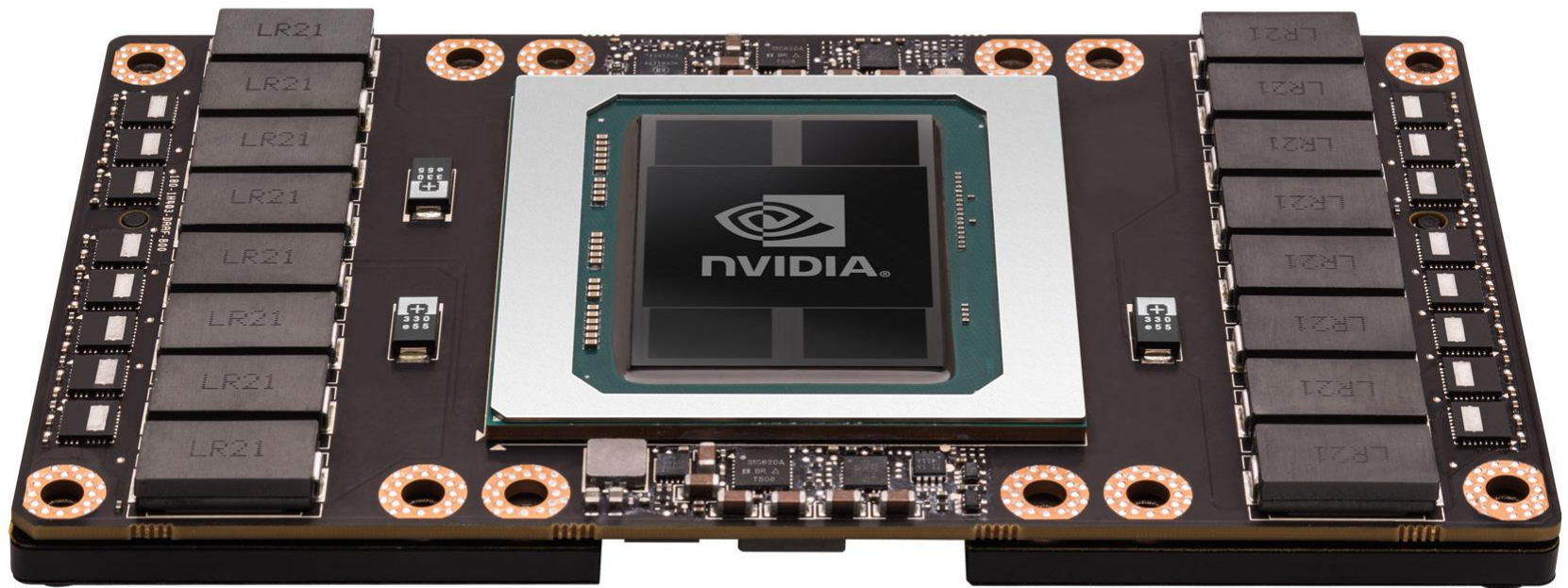
- Nvidia DGX-1 with panel open



http://www.anandtech.com/show/10229/nvidia-announces-dgx1-server

# GPU: System Organization

- Inside Nvidia DGX-1



7 TB SSD

8x Tesla P100 16GB

2x Xeon

3U - 3200W

NVLink Hybrid Cube Mesh

Quad IB 100Gbps, Dual 10GbE

# GPU: System Organization

- Tesla P100 in Nvidia DGX-1

# NVIDIA Tesla GP100



- 6 GPCs
- 30 TPCs
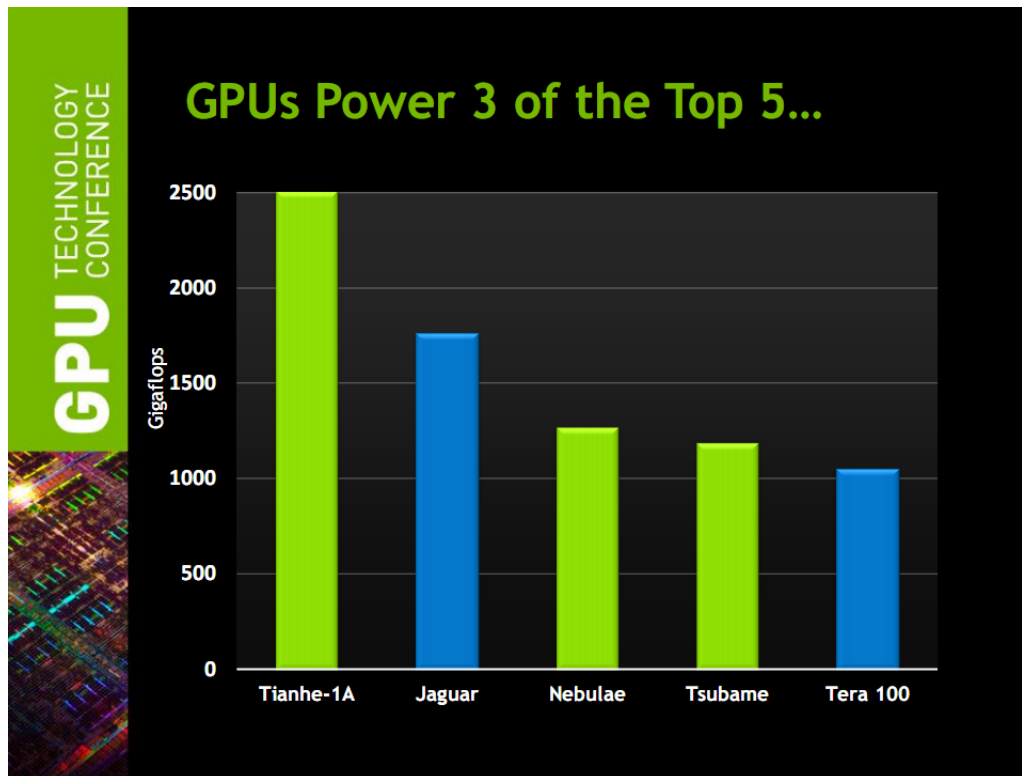- 60 SMs
- 3840 Cores
- 240 Tex
- 8 MCs
  - 4096-bit

6

# NVIDIA Telsa GP100 SM



- 64 FP32 cores
  - Supports 16-bit FP
- 32 FP64 cores
- 2 32K 32-bit Registers
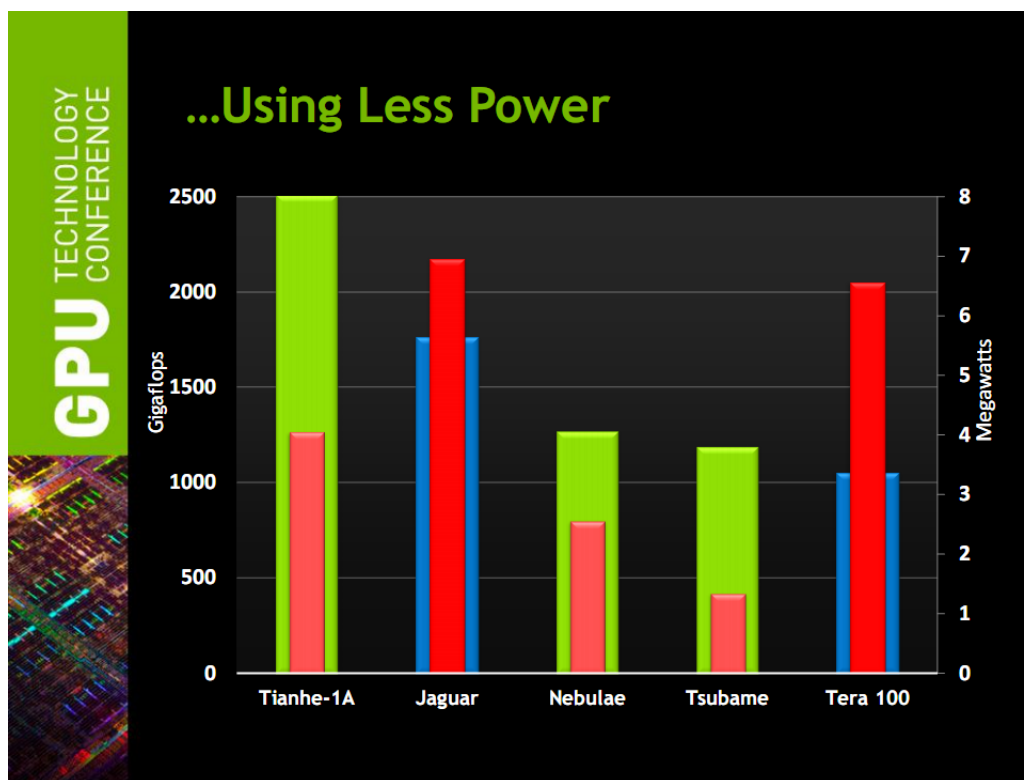- Separate shared mem.
- Unified Texture/L1 $

7

# GPUs Everywhere

- Many of the top-ranked supercomputers are based on GPU
  - The World #1 supercomputer is powered by NVIDIA GPU

# But Why are GPUs Everywhere?

- It is the performance/Watt metric
  - Mobile computers, Supercomputers, Laptops: Share power agony

# GPU: Initially for Graphics

- Take a scene and project on a screen
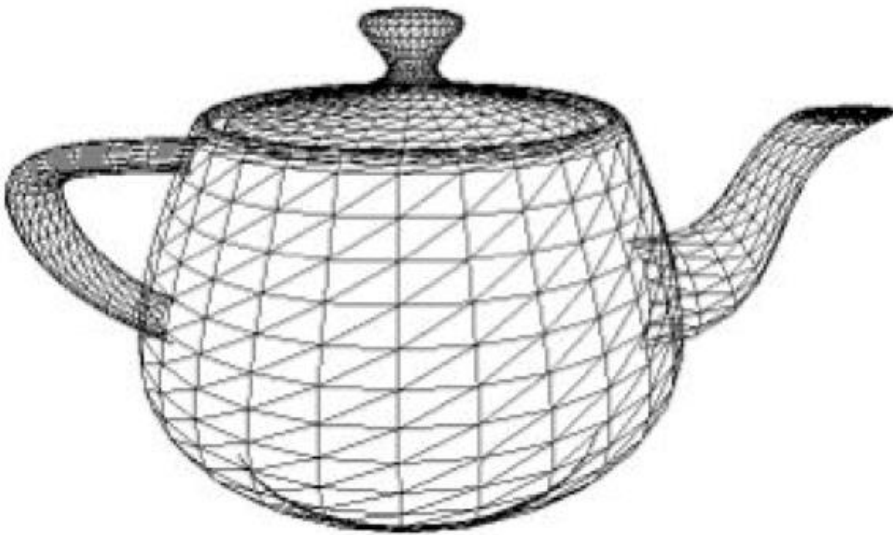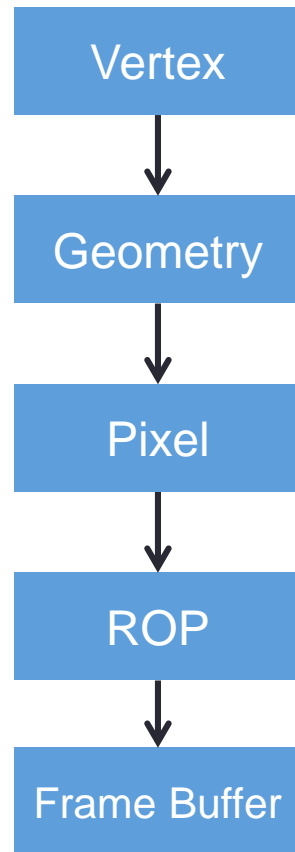  - How does each polygon translate into a screen pixel
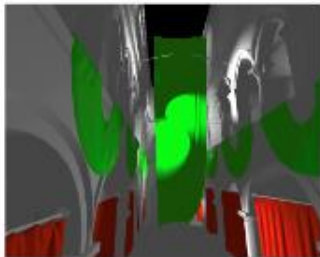


Image credit: Henrik Wann Jensen
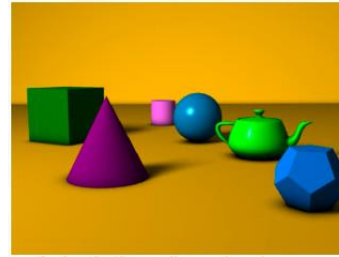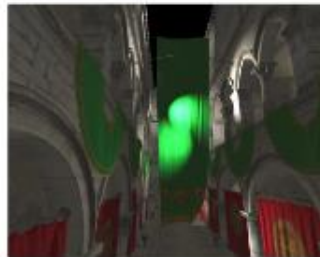
# Graphics-based GPU Pipeline (before 2006)



CS 550: Introduction to Computer Graphics

# GPU Pipeline Functions

- Vector Shader : Takes as input vectors given in 3D format (X,Y,Z) and projects to a 2D space

- Pixel Shader: Takes each projected point in 2D space and puts texture, color and depth for each point
  - Texture mapping: Adding an image on top of a scene
  - Depth mapping: Shows which objects will obscure other objects



Example Image: Before and After Texture Mapping

Example Image: Z-Buffer Representation of Image

- Geometry Shaders: Operates directly on graphics primitives, such as lines, triangles (as opposed to pixels)

CS 550: Introduction to Computer Graphics

# Need for Unified Shader



GPU Architecture: Implications & Trends D. Luebke, Nvidia
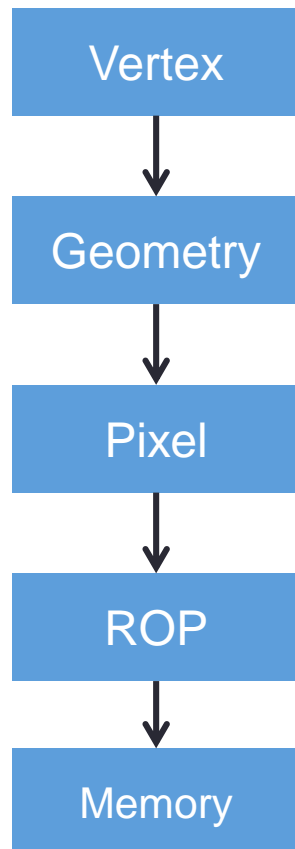
# Resource Balancing Unified Shader
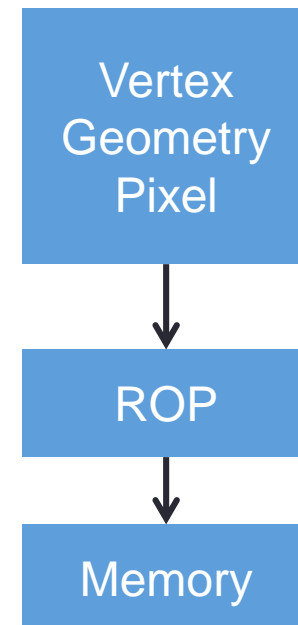
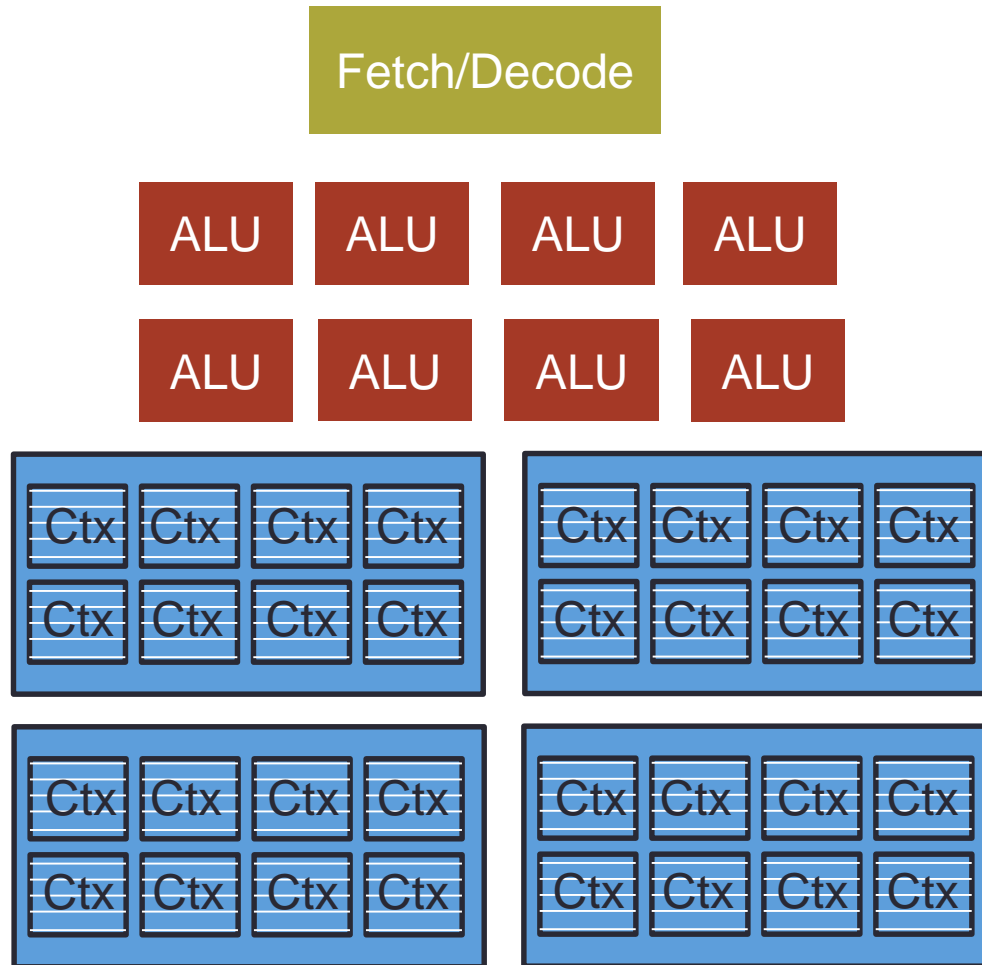# GPU Pipeline evolution
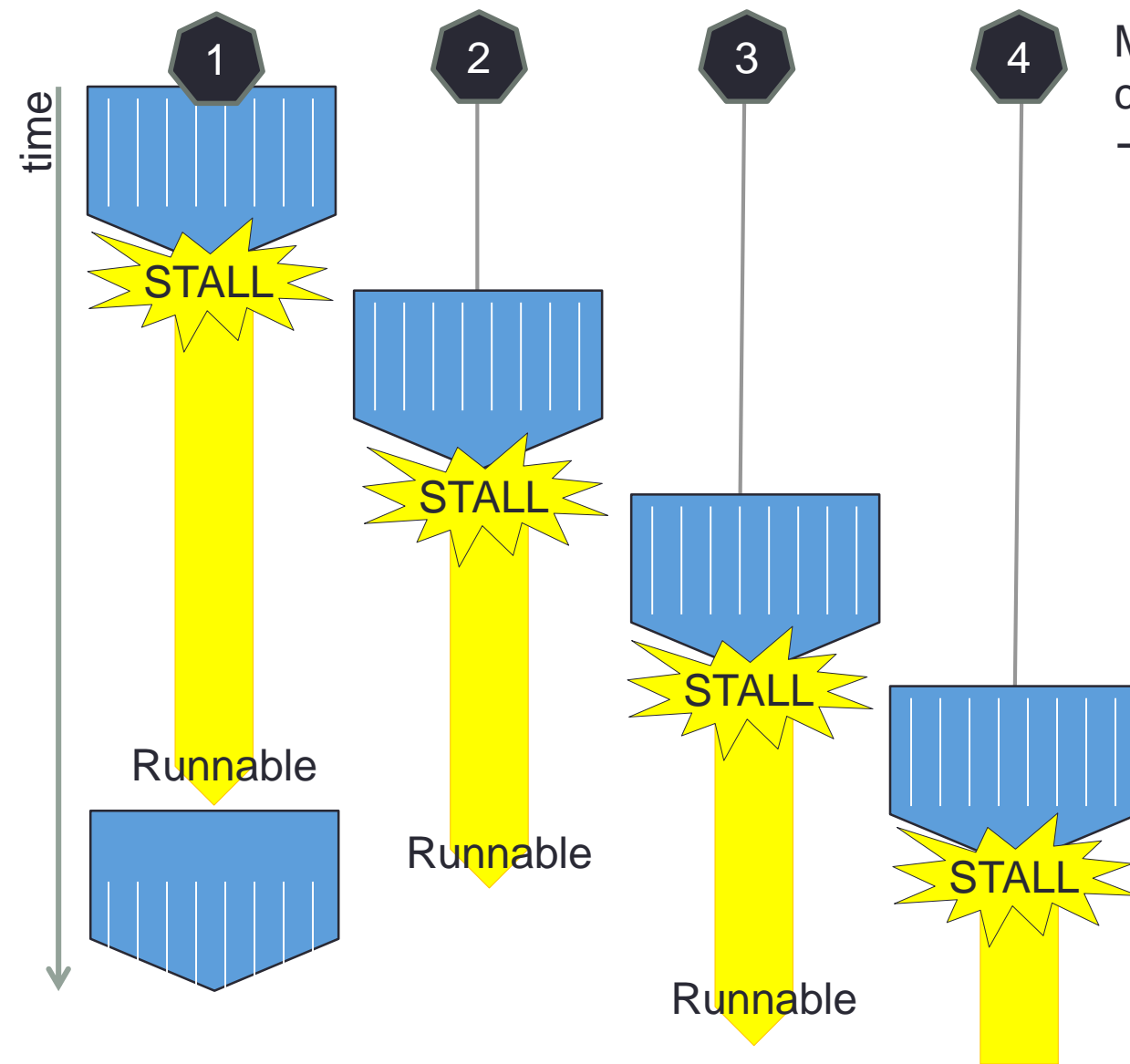
**Before 2006**

- Processor per function

**After 2006**

- Unified Shader Model(GeForce 8 Series)
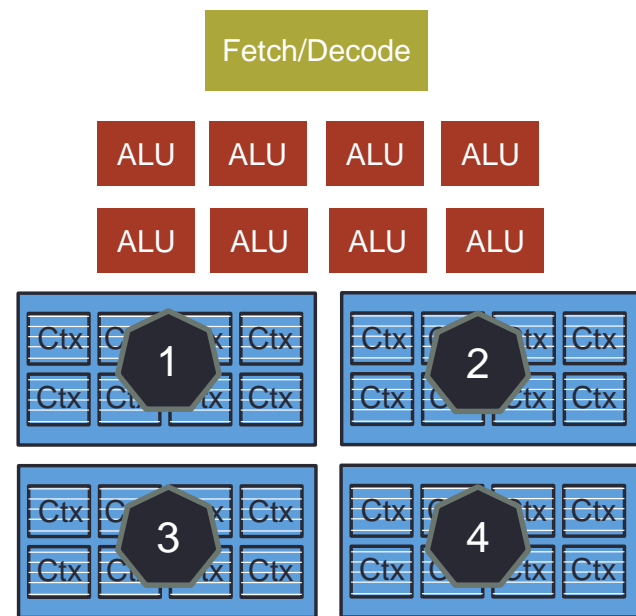- Any work can be performed on any shader core → High performance computing
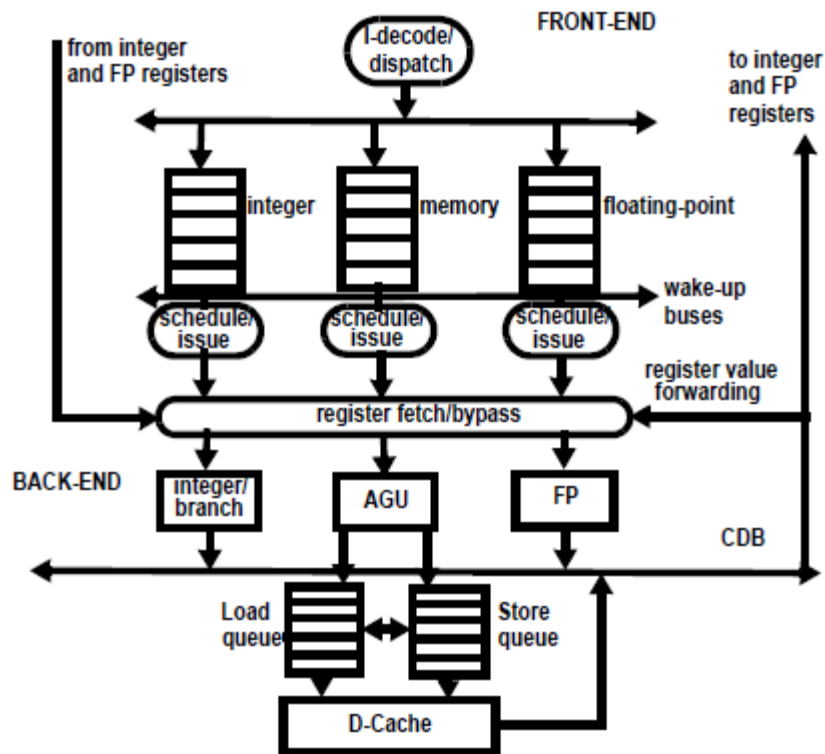
```
Vertex
  ↓
Geometry
  ↓
Pixel
  ↓
ROP
  ↓
Memory
```

```
Vertex
Geometry
Pixel
  ↓
ROP
  ↓
Memory
```

15

# Unified Shader Model

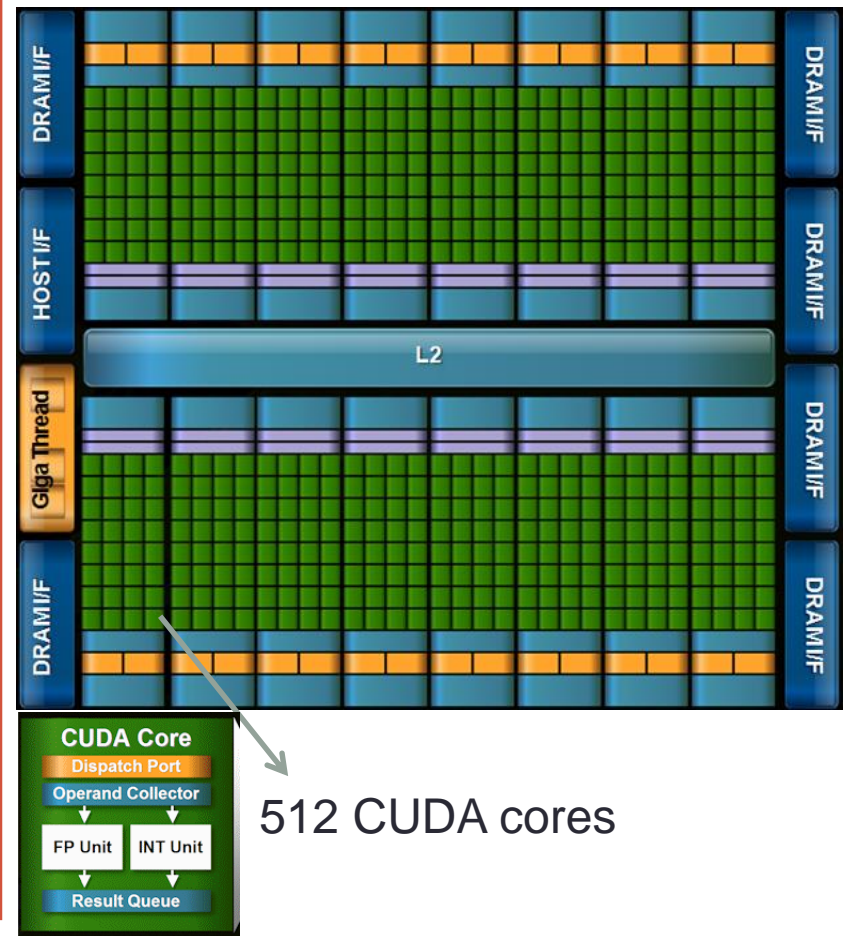# Unified Shader Model

# GPU vs CPU

## Complex OoO CPU



## GPU (Nvidia FERMI)



512 CUDA cores

# Basic idea of high performance GPU

- Many simple cores
  - No fancy Branch prediction
  - No Complex O-o-O control logic
  - No memory prefetcher
  - No cache coherence
  - . . .
- Unified shader cores
  - Sharing instruction stream across groups of fragments
- Stall latency hiding
  - When a group stalls, work on another group