



CSCI-GA.3033-012

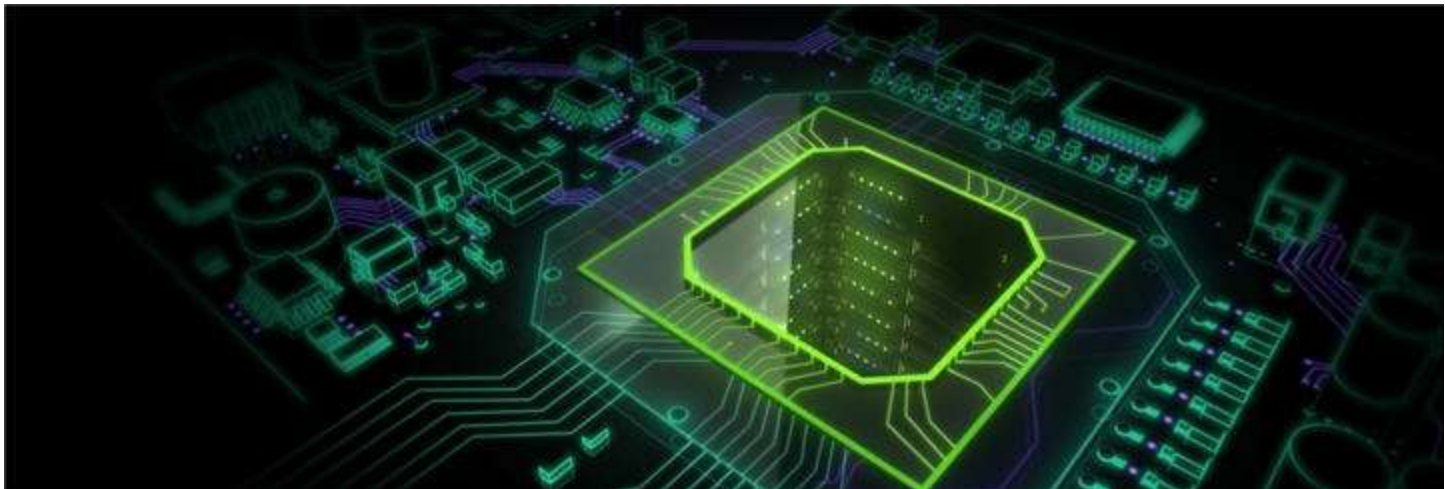
# Graphics Processing Units (GPUs): Architecture and Programming

## Lecture 2: History of GPU Computing

Mohamed Zahran (aka Z)

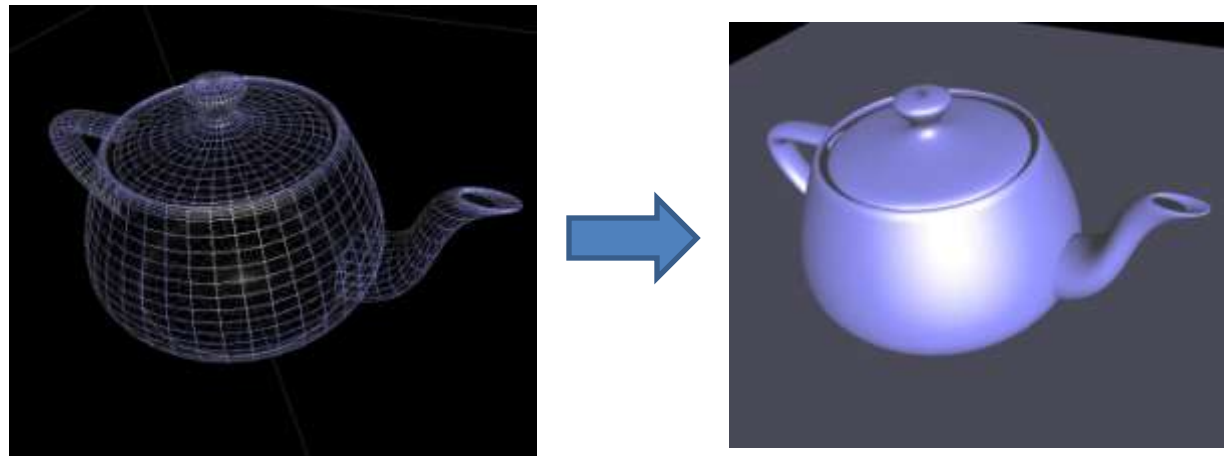
mzahran@cs.nyu.edu

<http://www.mzahran.com>



# A Little Bit of Vocabulary

- **Rendering**: the process of generating an image from a model
- **Vertex**: the corner of a polygon (usually that polygon is a triangle)
- **Pixel**: smallest addressable screen element



# From Numbers to Screen

```
0.748952 -0.764952 -0.210132,  
0.872246 -0.600062 -0.210132,  
1.00016 -0.369596 -0.210132,  
1.00939 -0.004004 -0.210132,  
1.14496 0.324436 -0.210132,  
1.15747 0.601712 -0.210132,  
1.03010 0.793529 -0.210132,  
0.09164 0.972032 -0.210132,  
0.808203 0.929010 -0.210132,  
0.442563 0.985585 -0.210132,  
0.221794 1.00159 -0.210132,  
0 1.0053 -0.210132,  
-0.221794 1.00159 -0.210132,  
-0.442563 0.985585 -0.210132,  
-0.808203 0.929010 -0.210132,
```

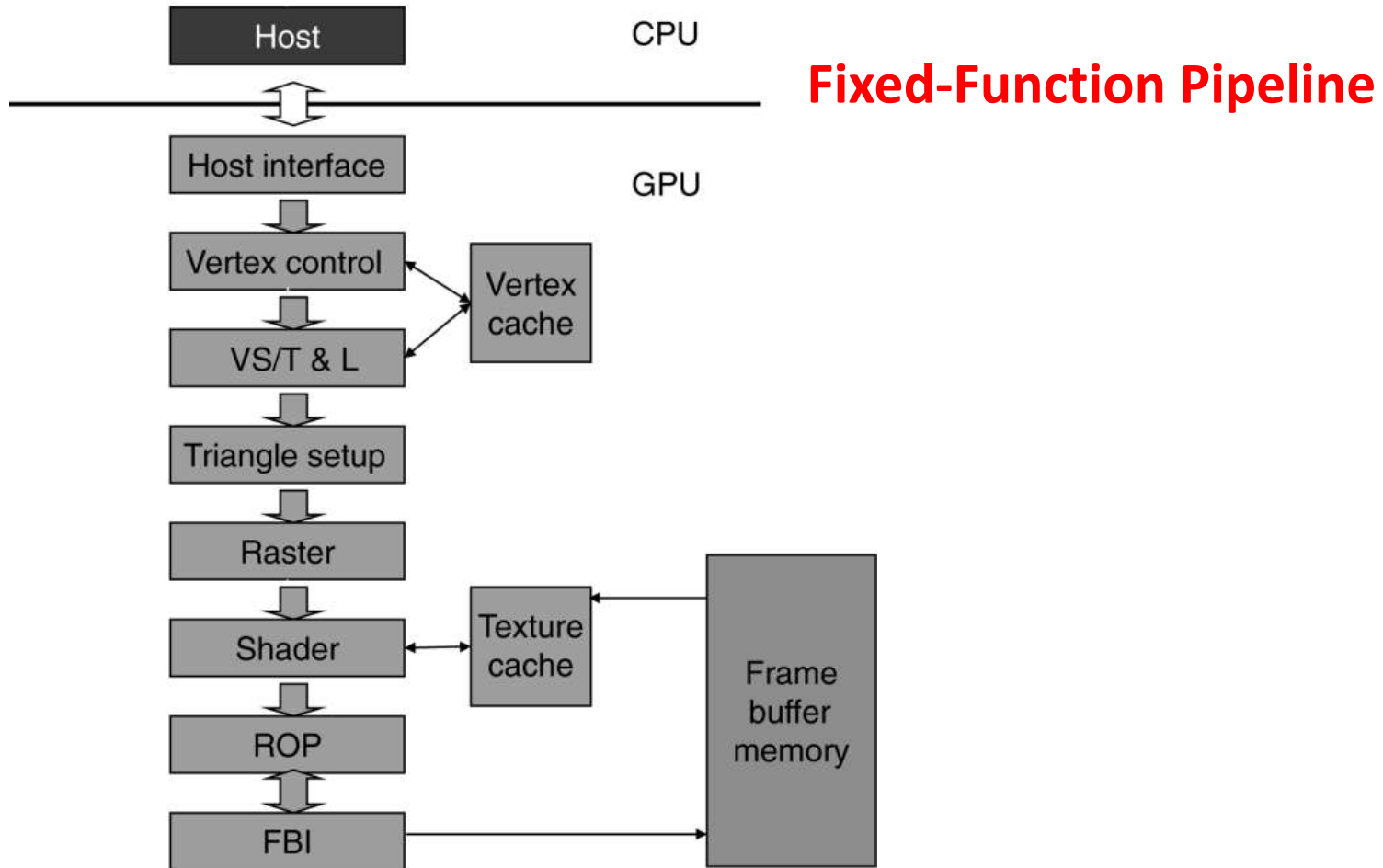


# Before GPUs

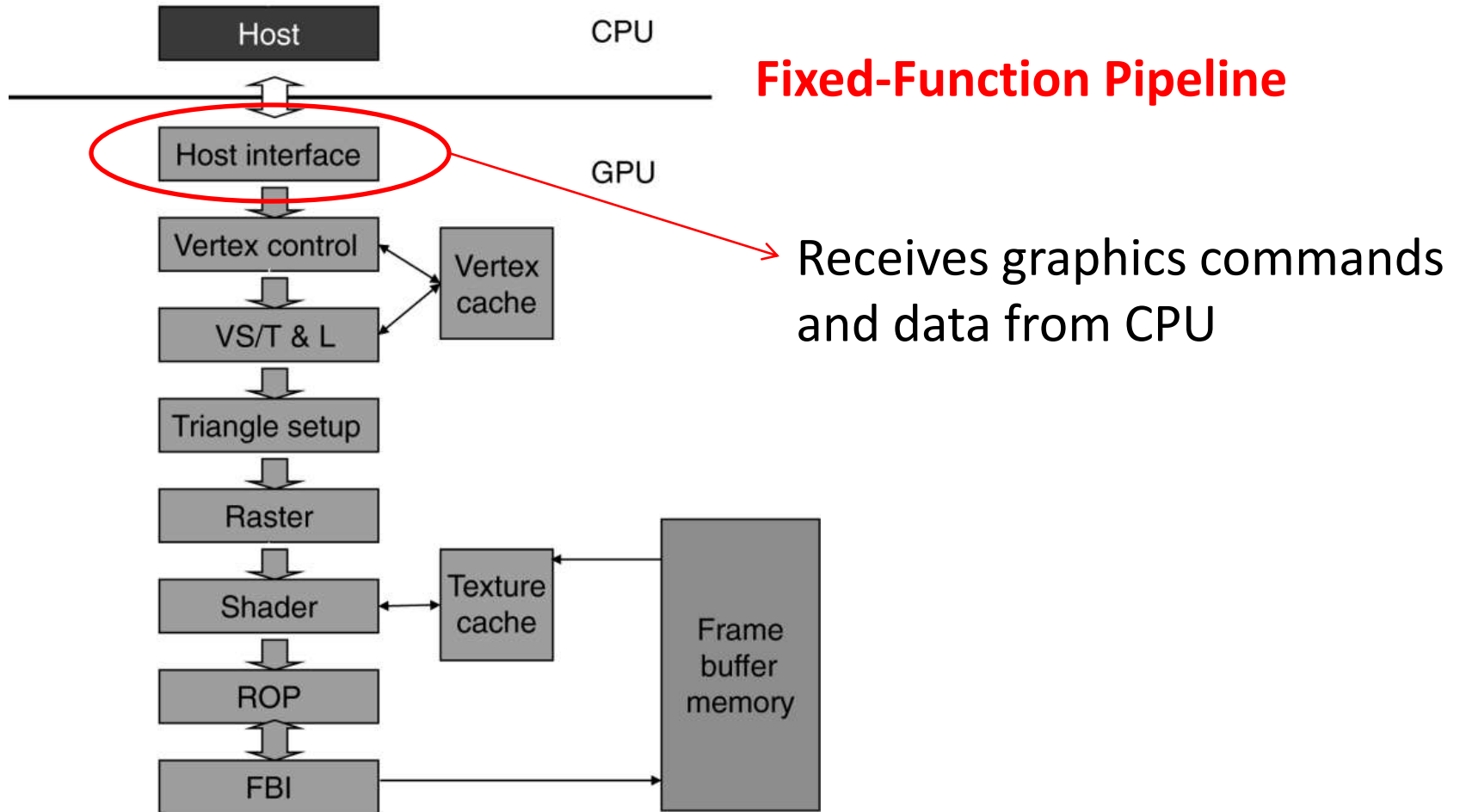
- Vertices to pixels:
  - Transformations done on CPU
  - Compute each pixel "by hand", in series...  
**slow!**

**Example:** 1 million triangles \* 100 pixels  
per triangle \* 10 lights \* 4 cycles per  
light computation = **4 billion cycles**

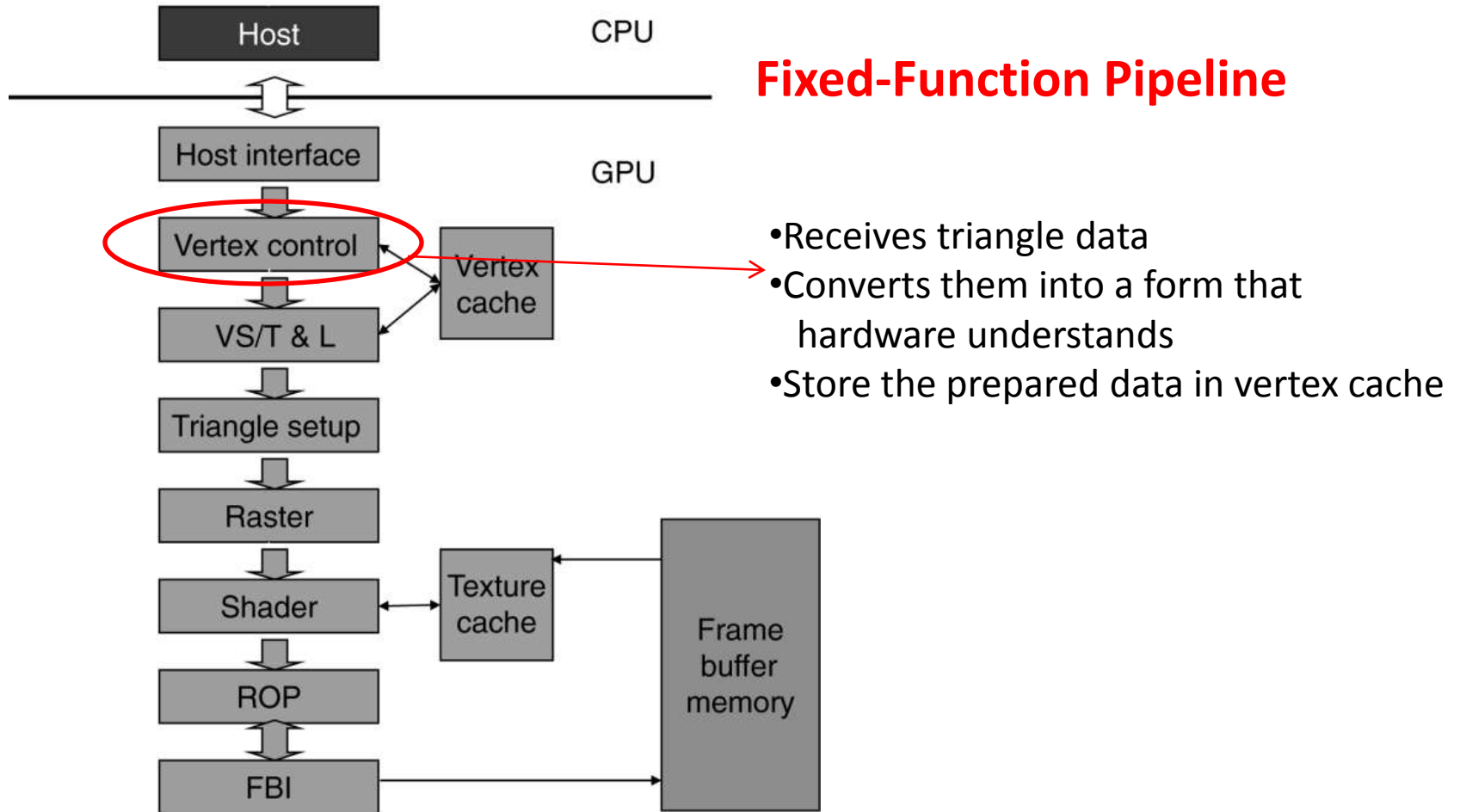
# Early GPUs: Early 80s to Late 90s



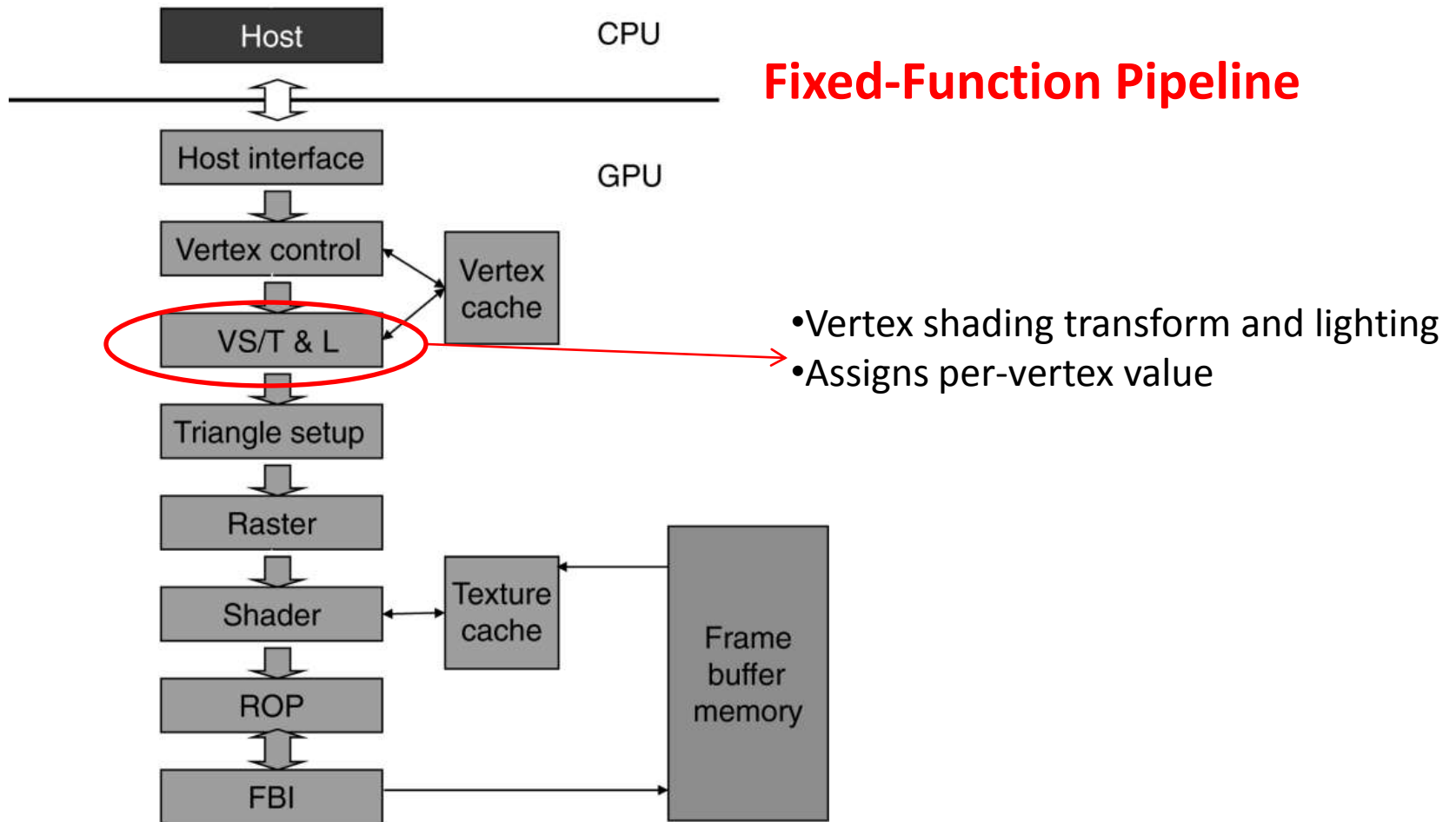
# Early GPUs: Early 80s to Late 90s



# Early GPUs: Early 80s to Late 90s

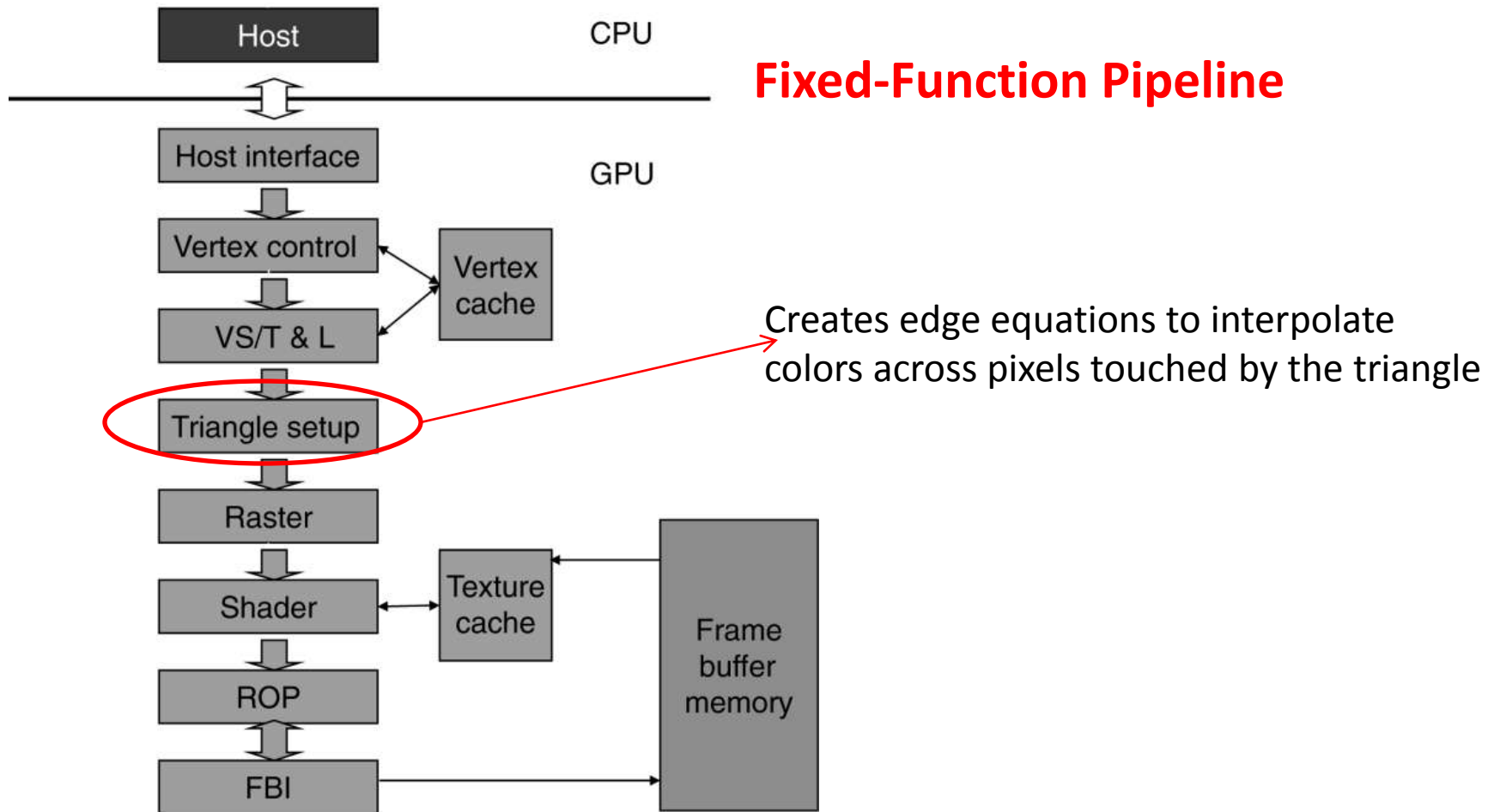


# Early GPUs: Early 80s to Late 90s

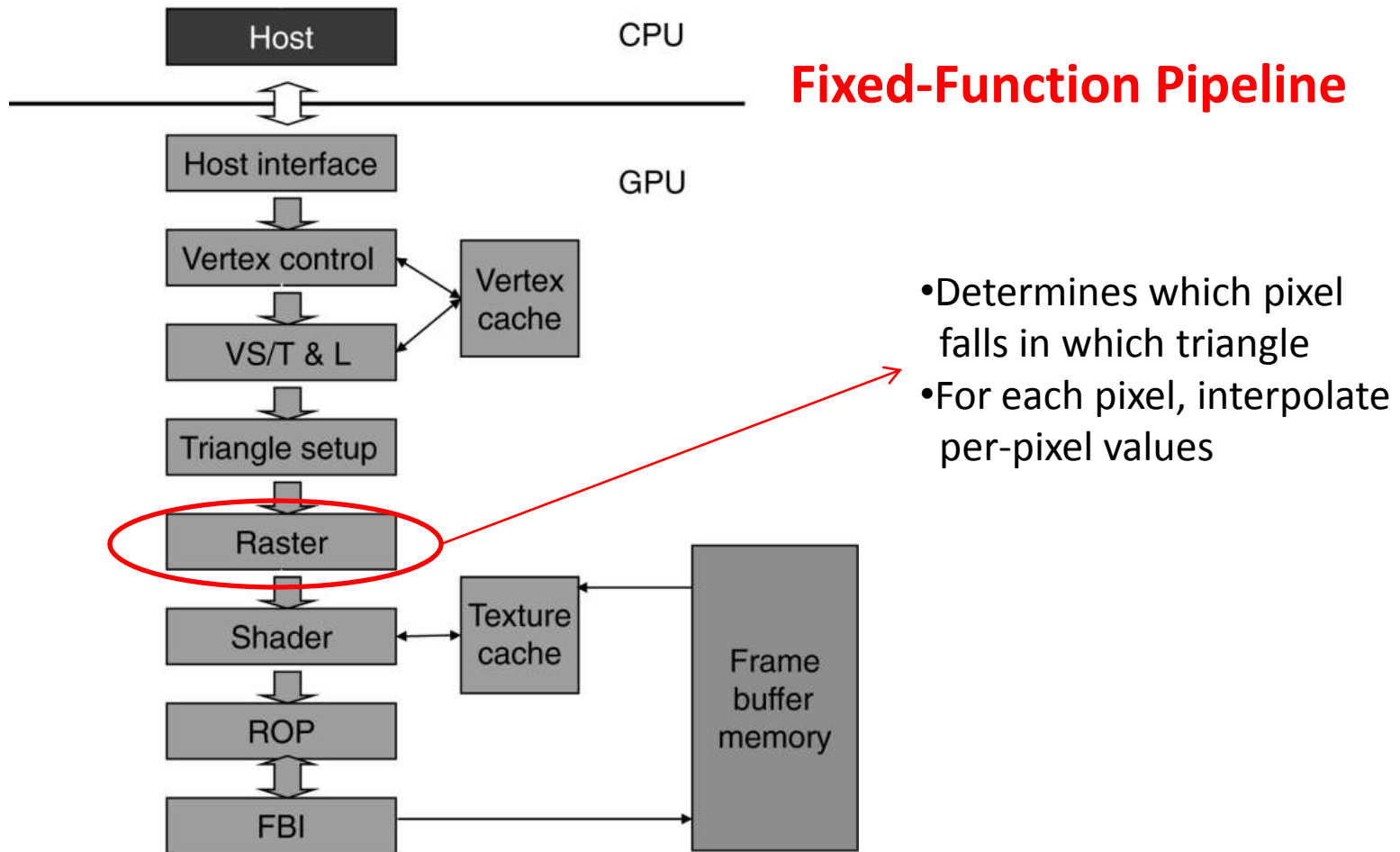




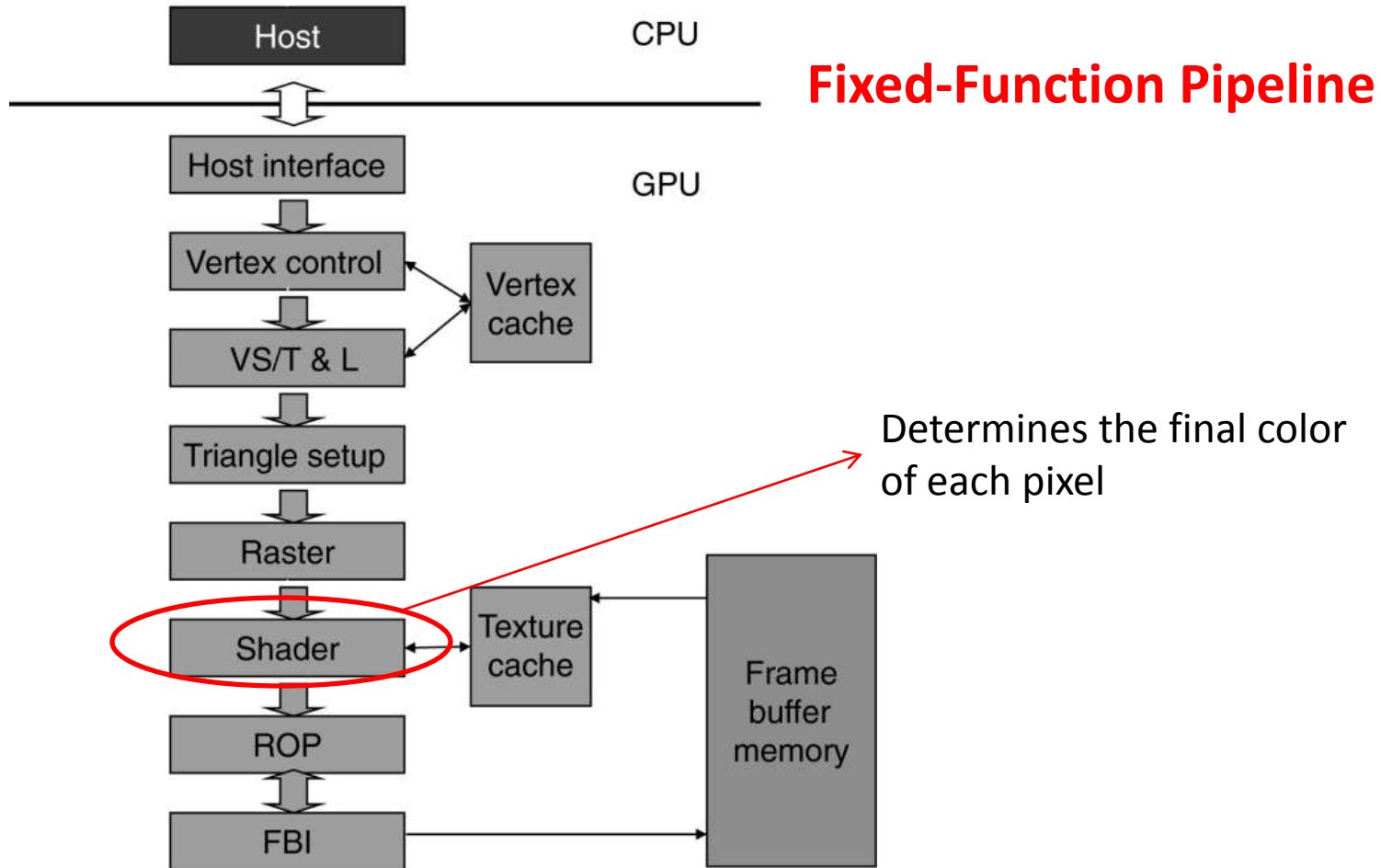
# Early GPUs: Early 80s to Late 90s



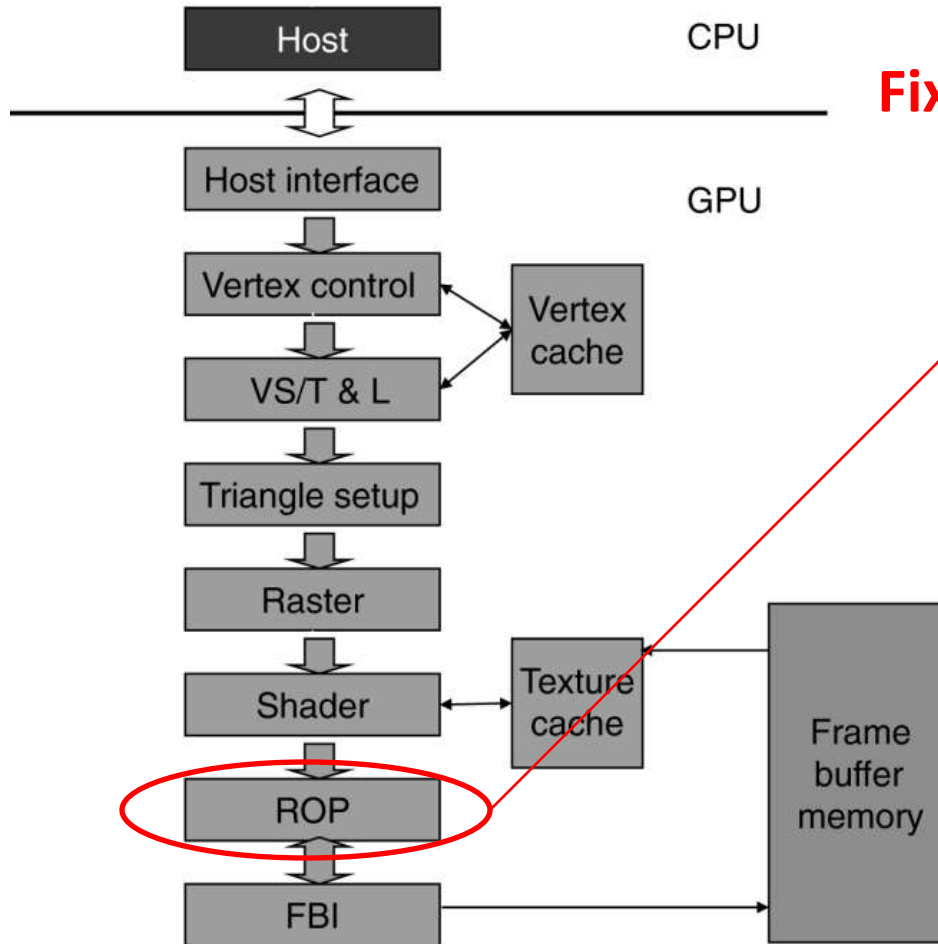
# Early GPUs: Early 80s to Late 90s



# Early GPUs: Early 80s to Late 90s

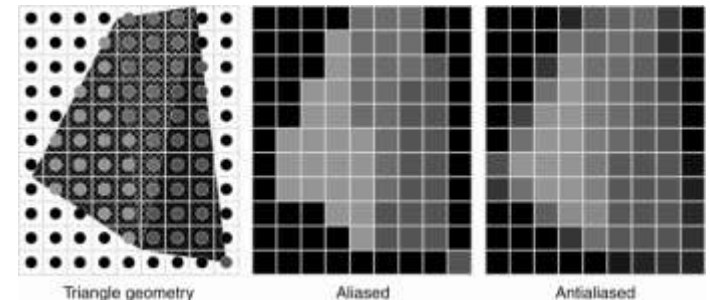


# Early GPUs: Early 80s to Late 90s

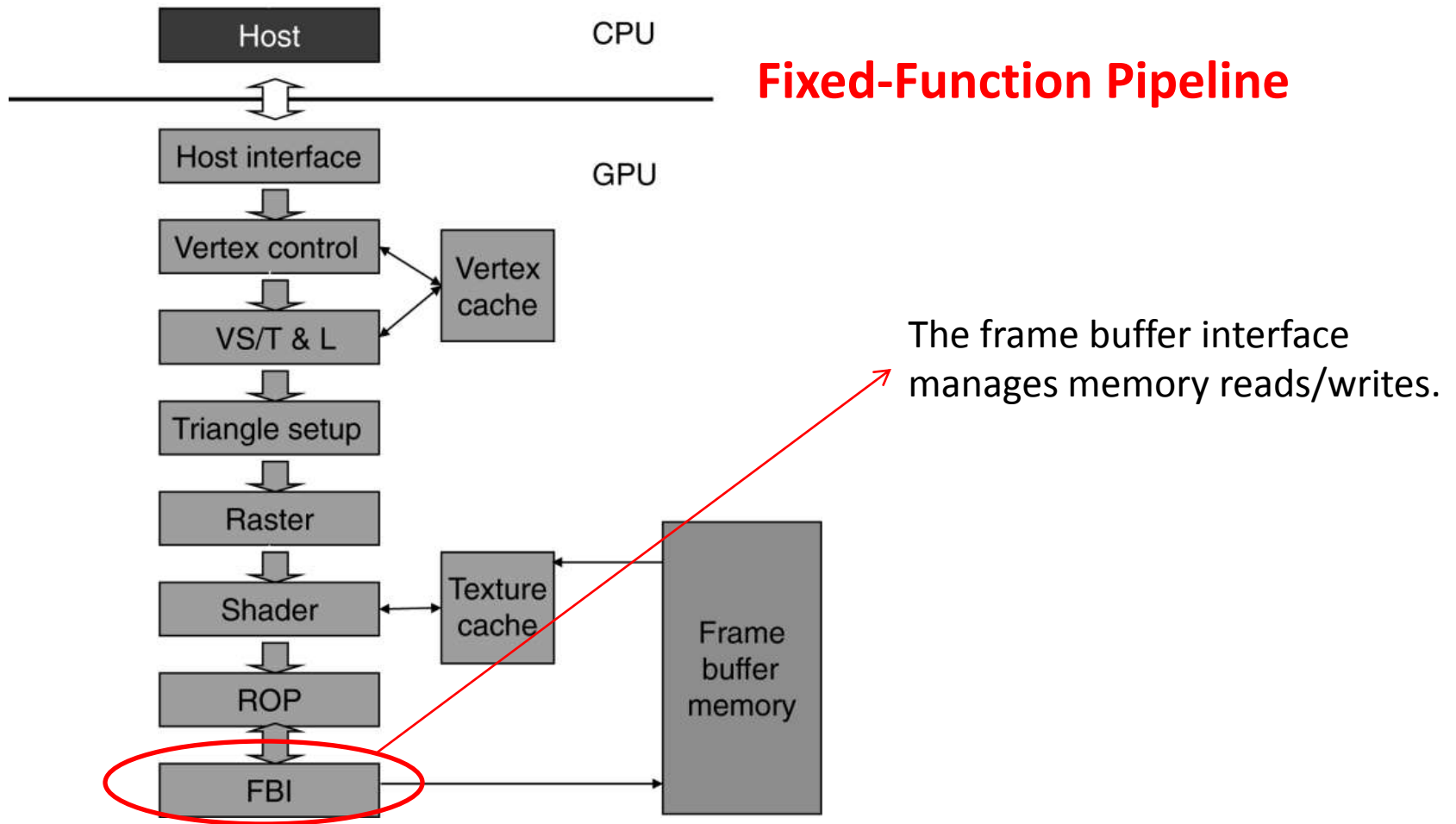


## Fixed-Function Pipeline

The raster operation:  
performs color raster operations  
that blend the color of overlapping  
objects for transparency and  
antialiasing

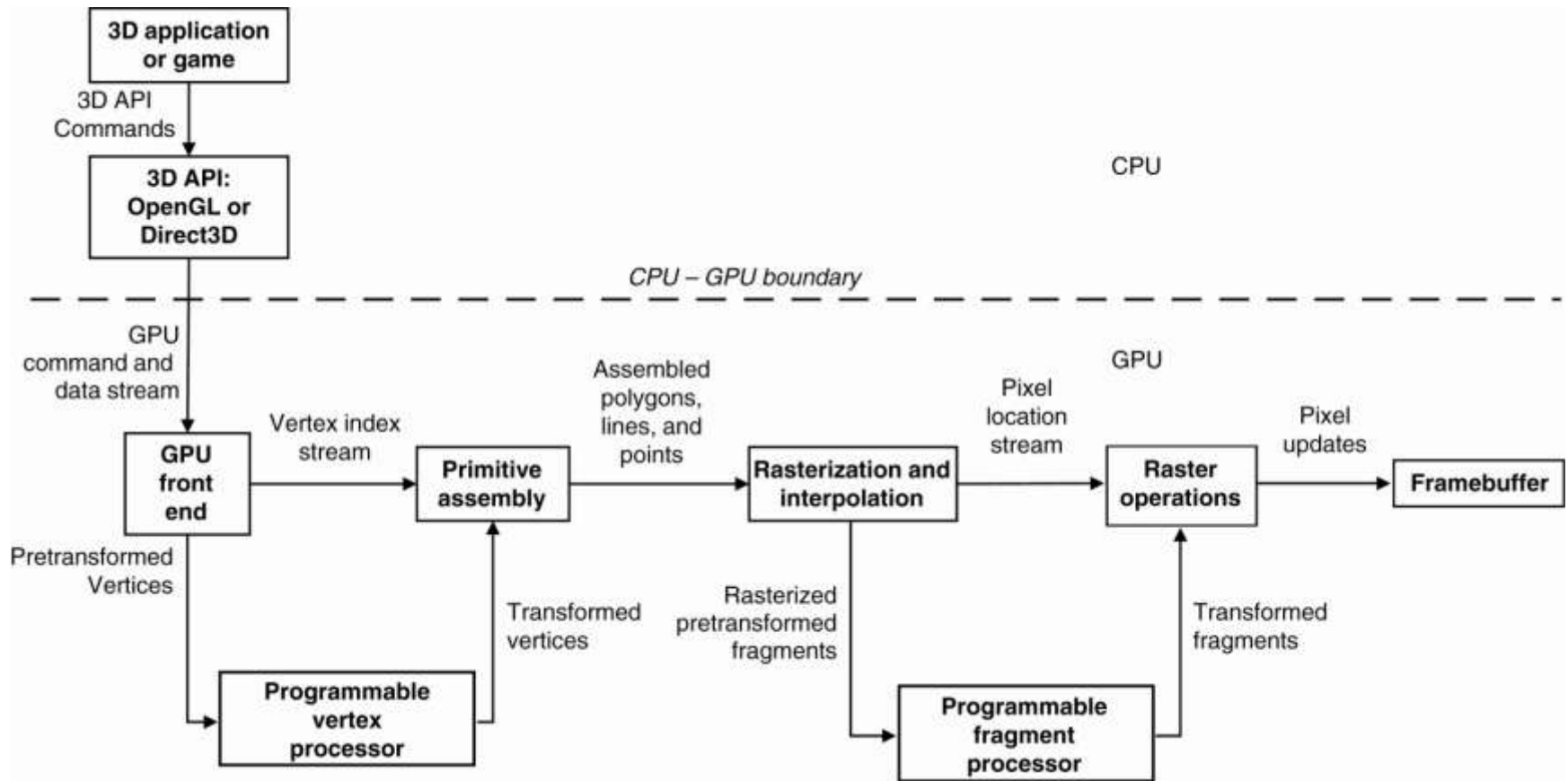


# Early GPUs: Early 80s to Late 90s



# Next Steps

- In 2001:
  - NVIDIA exposed the application developer to the instruction set of VS/T&L stage
- Later:
  - General programmability extended to to shader stage
  - Data independence is exploited



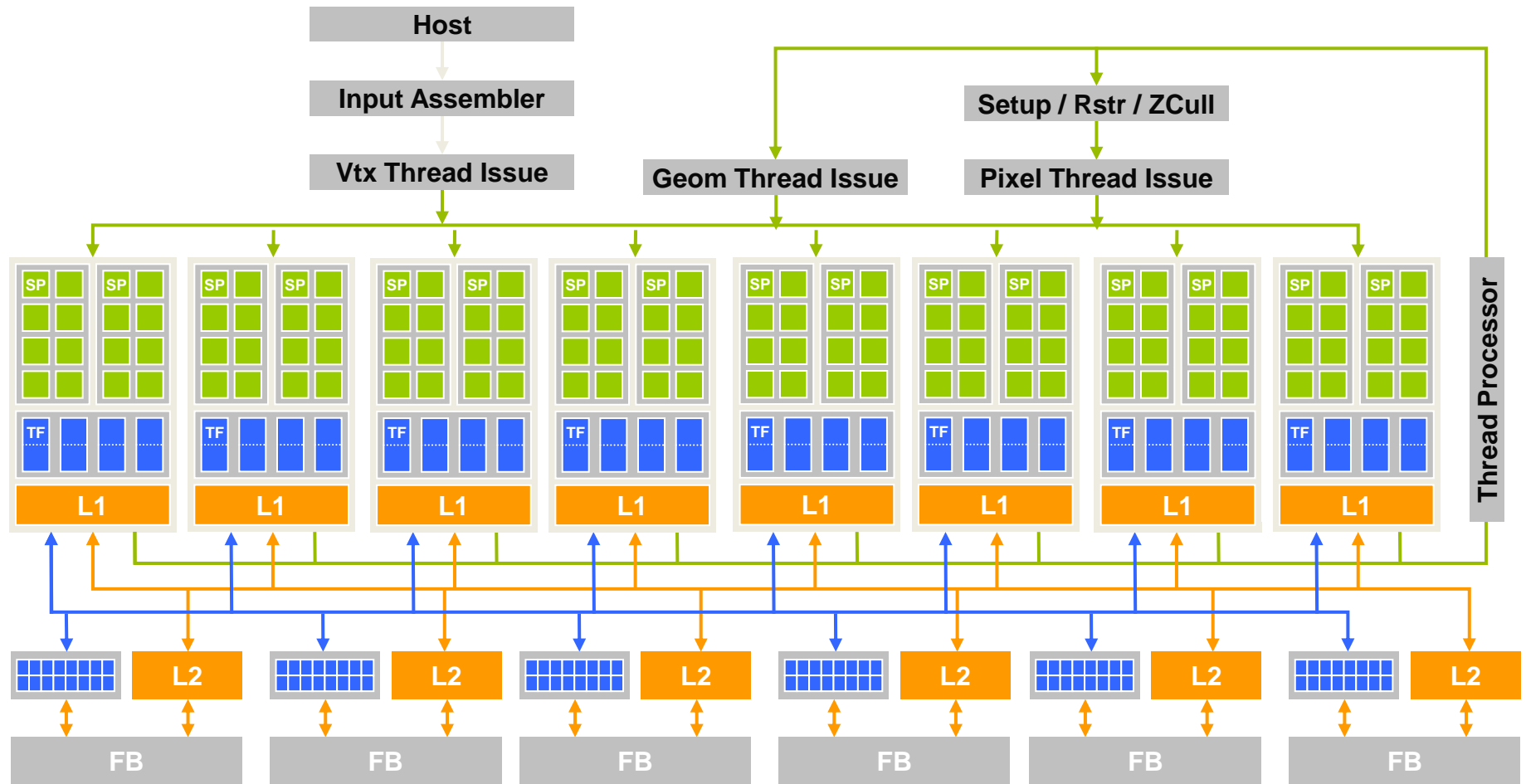
# In 2006

- NVIDIA GeForce 8800 mapped separate graphics stage to a **unified array of processors**
  - For vertex shading, geometry processing, and pixel processing
  - Allows dynamic partition



# Regularity + Massive Parallelism





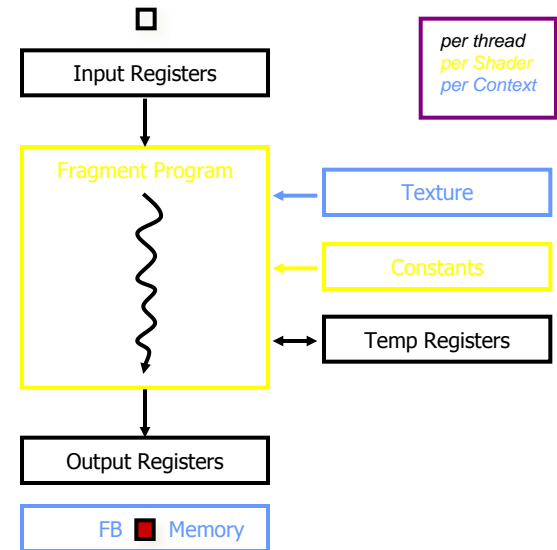
Exploring the use of GPUs to solve compute intensive problems

The birth of GPGPU but there are many constraints

GPUs and associated APIs were designed to process graphics data

# Previous GPGPU Constraints

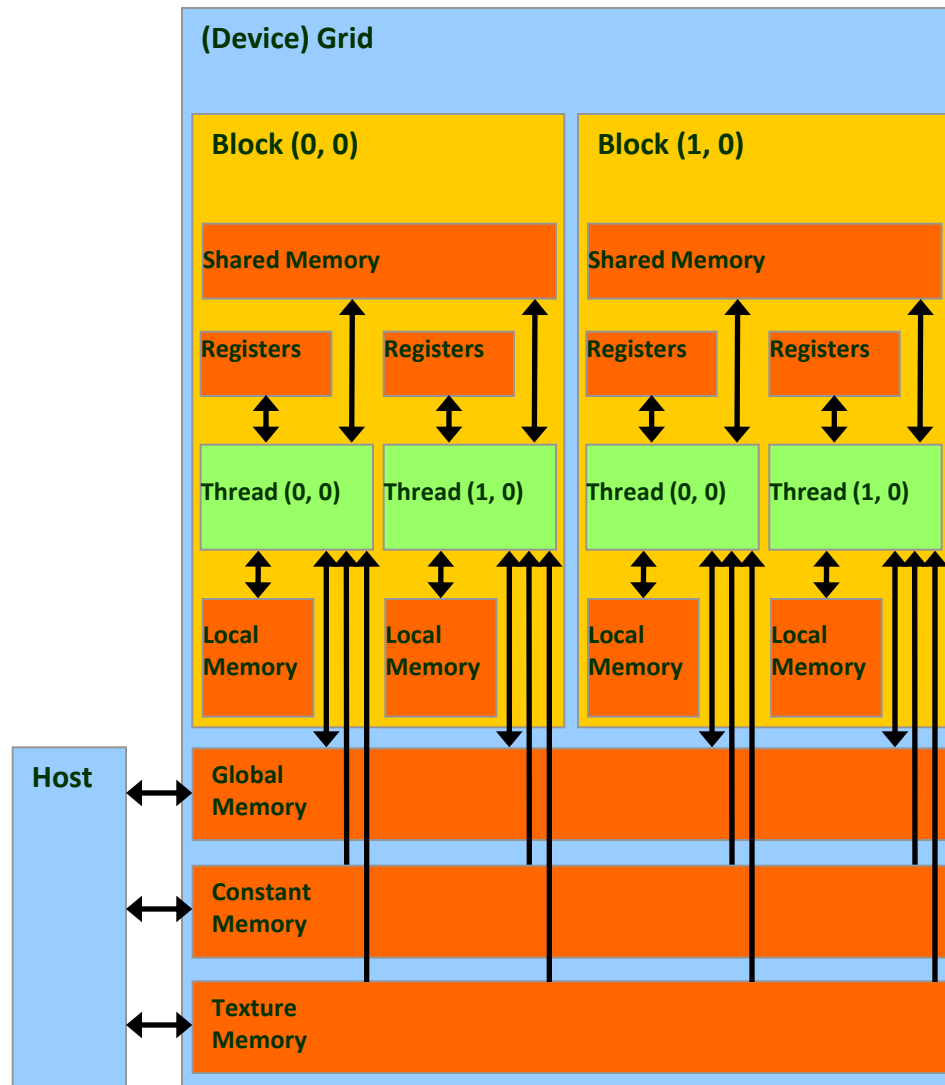
- Dealing with graphics API
  - Working with the corner cases of the graphics API
- Addressing modes
  - Limited texture size/dimension
- Shader capabilities
  - Limited outputs
- Instruction sets
  - Lack of Integer & bit ops
- Communication limited
  - Between pixels
  - Scatter  $a[i] = p$



# The Birth of GPU Computing

- **Step 1:** Designing high-efficiency floating-point and integer processors.
- **Step 2:** Exploiting data parallelism by having large number of processors
- **Step 3:** Shader processors fully programmable with large instruction cache, instruction memory, and instruction control logic.
- **Step 4:** Reducing the cost of hardware by having multiple shader processors to share their cache and control logic.
- **Step 5:** Adding memory load/store instructions with random byte addressing capability
- **Step 6:** Developing CUDA C/C++ compiler, libraries, and runtime software models.

# A Glimpse on Memory Space

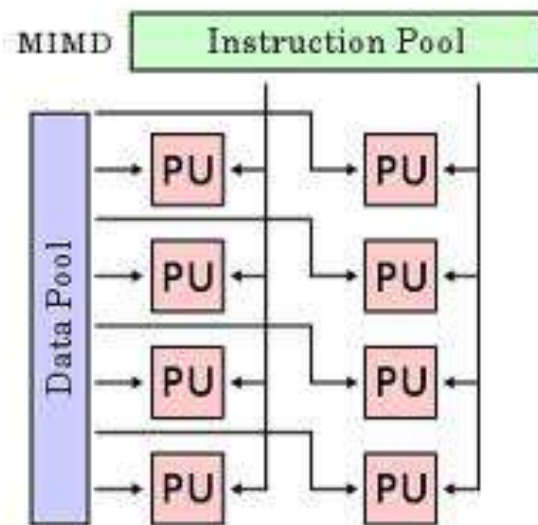
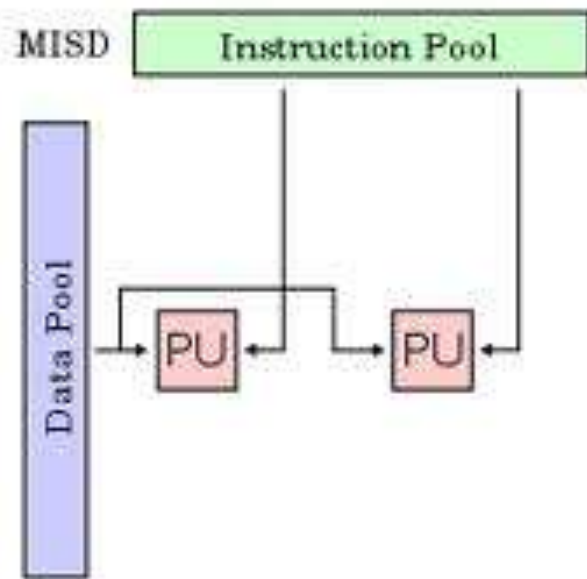
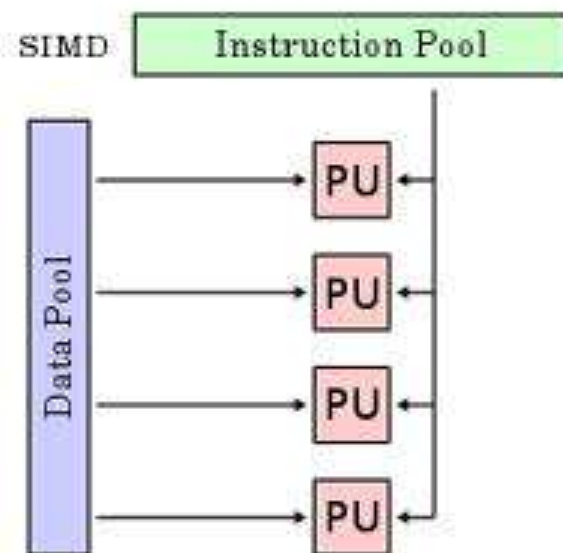
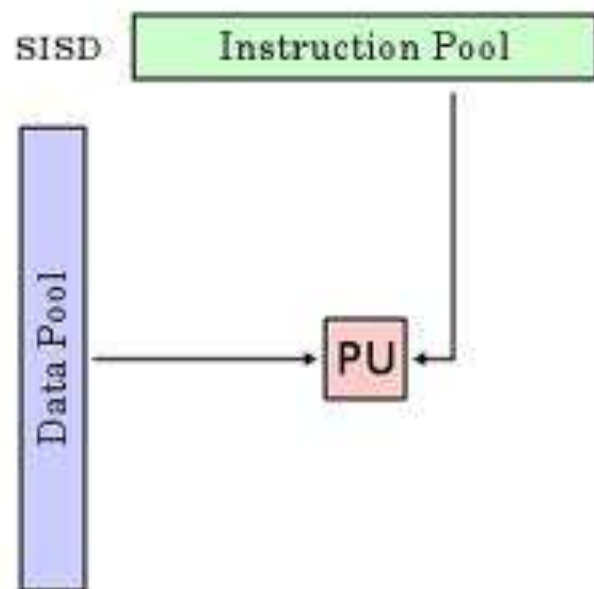


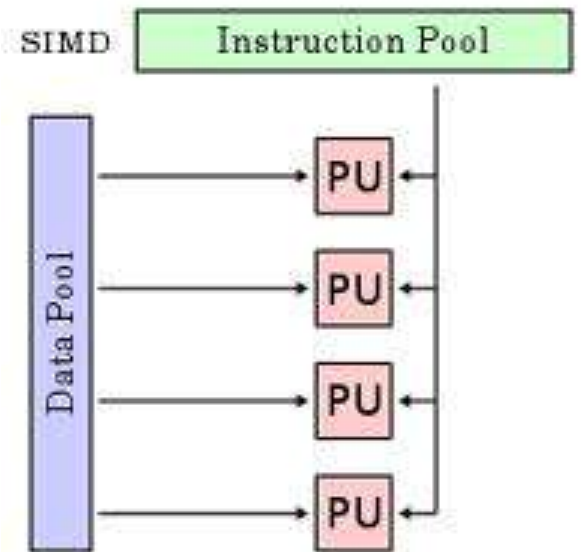
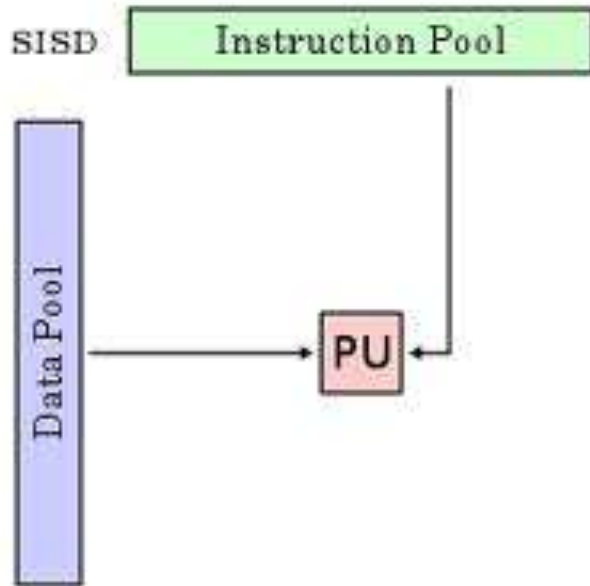
Source : "NVIDIA CUDA Programming Guide" version 1.1

# A Quick Glimpse on: Flynn Classification

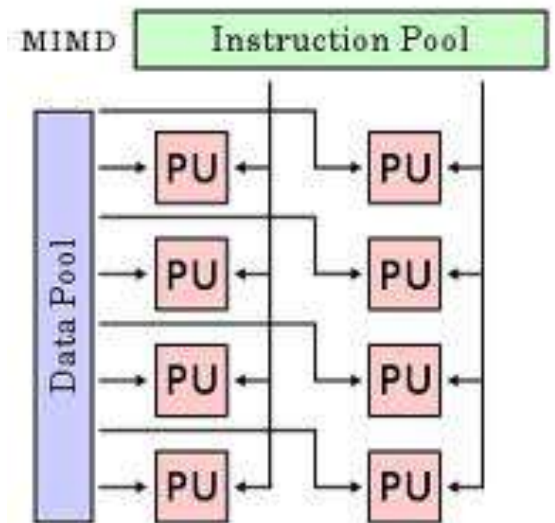
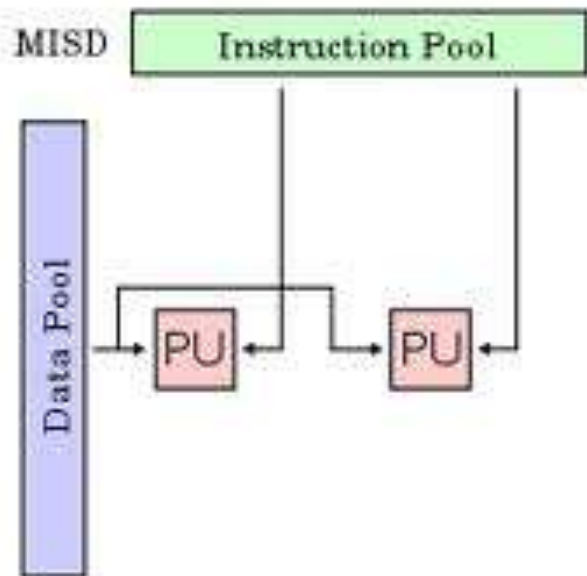
- A taxonomy of computer architecture
- Proposed by Micheal Flynn in 1966
- It is based two things:
  - Instructions
  - Data

	Single instruction	Multiple instruction
Single data	<b>SISD</b>	<b>MISD</b>
Multiple data	<b>SIMD</b>	<b>MIMD</b>



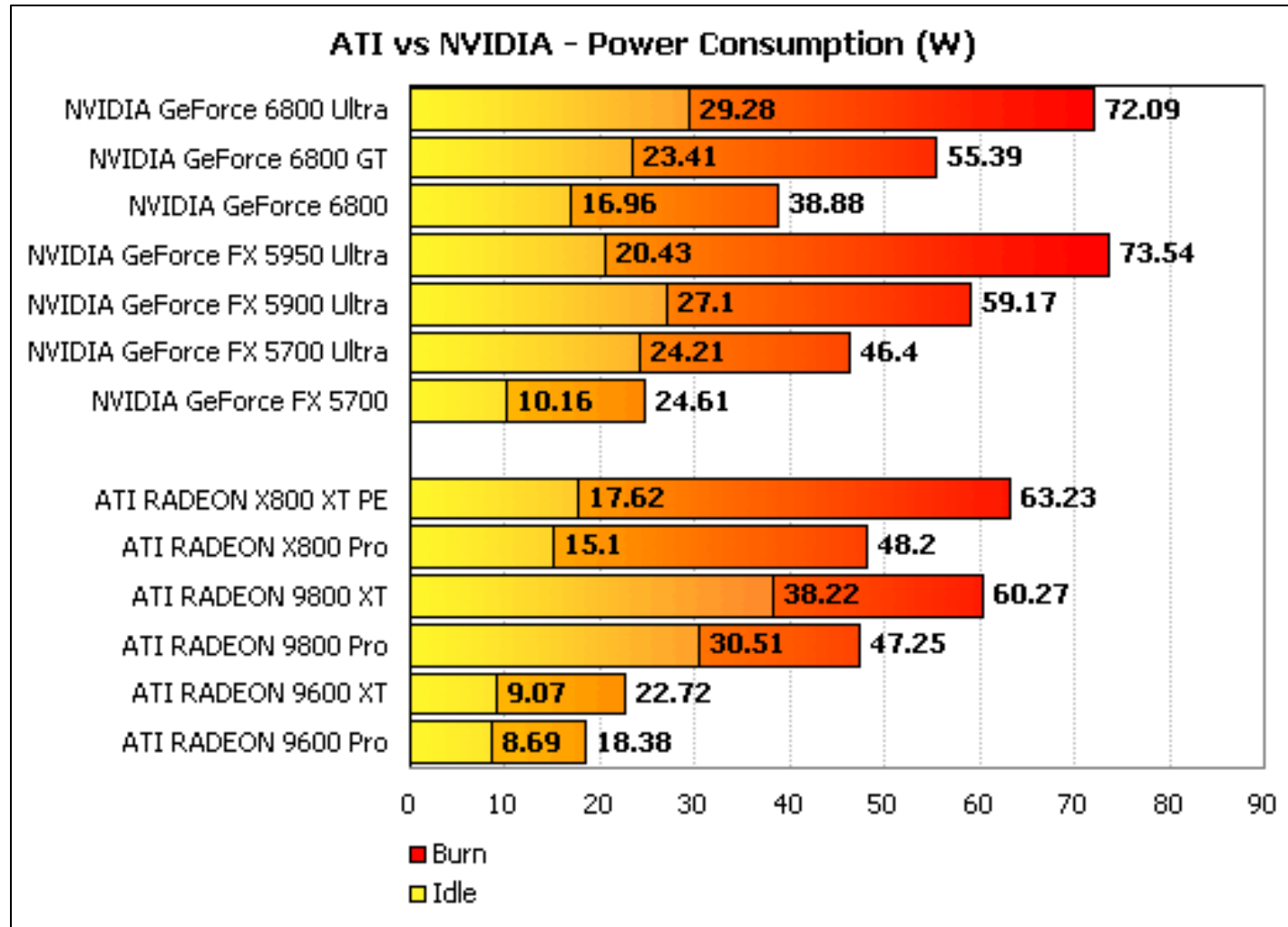


Which one  
is closest to  
GPU?





# Problem With GPUs: Power



# Problems with GPUs

- Need enough parallelism
- Under-utilization
- Bandwidth to CPU

Still a way to go

# Conclusions

- The design of state-of-the art GPUs includes:
  - Data parallelism
  - Programmability
  - Much less restrictive instruction set