# GPGPU introduction and network applications

PacketShaders, SSLShader

**Mellanox** TECHNOLOGIES

Connect. Accelerate. Outperform.™

# Agenda

- **GPGPU Introduction**
  - Computer graphics background
  - GPGPUs – past, present and future
- **PacketShader – A GPU-Accelerated Software Router**
- **SSLShader – A GPU-Accelerated SSL encryption/decryption proxy**

# GPGPU Intro

# GPU = Graphics Processing Unit

- The heart of graphics cards

- Mainly used for real-time 3D game rendering
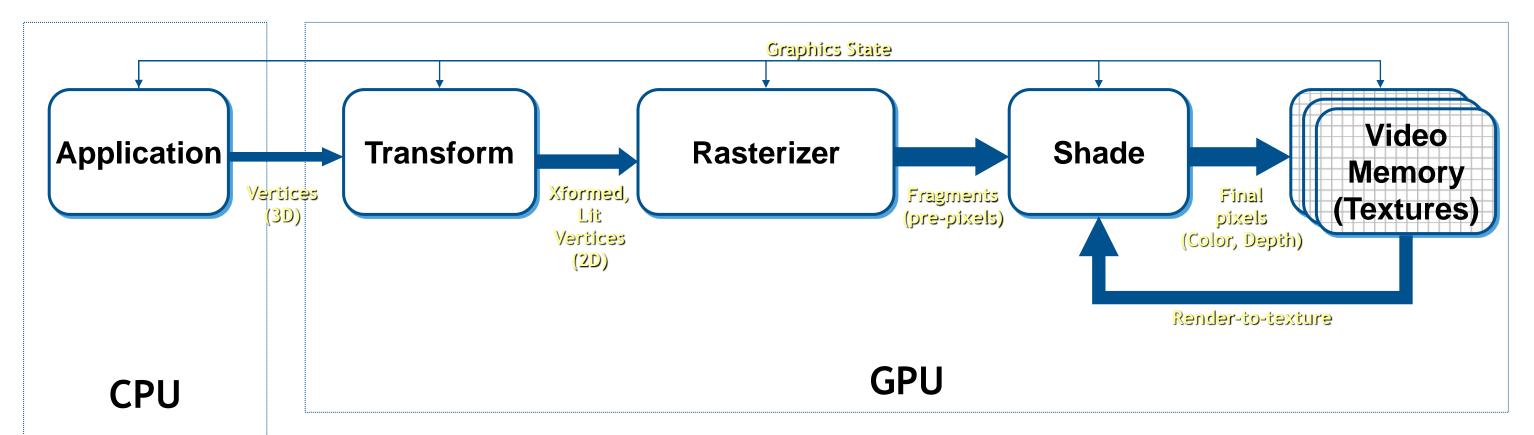  - Massively-parallel processing capacity



**(Ubisoft's AVARTAR, from http://ubi.com)**

# GPU Fundamentals: The Graphics Pipeline

Graphics State

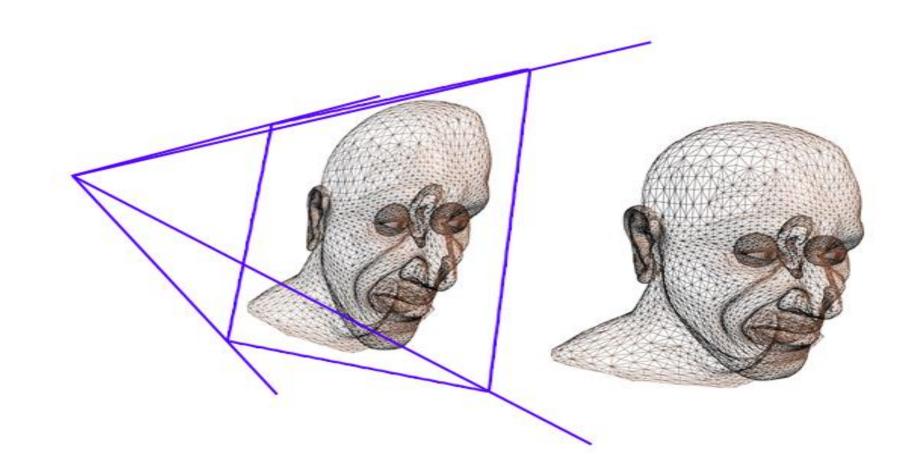| Application | → Vertices (3D) → | Transform | → Xformed, Lit Vertices (2D) → | Rasterizer | → Fragments (pre-pixels) → | Shade | → Final pixels (Color, Depth) → | Video Memory (Textures) |

Render-to-texture

**CPU**

**GPU**

- A simplified graphics pipeline
  - Note that pipe widths vary
  - Many caches, FIFOs, and so on not shown

- Vertex Processor (multiple operate in parallel)
  - Transform from "world space" to "image space"
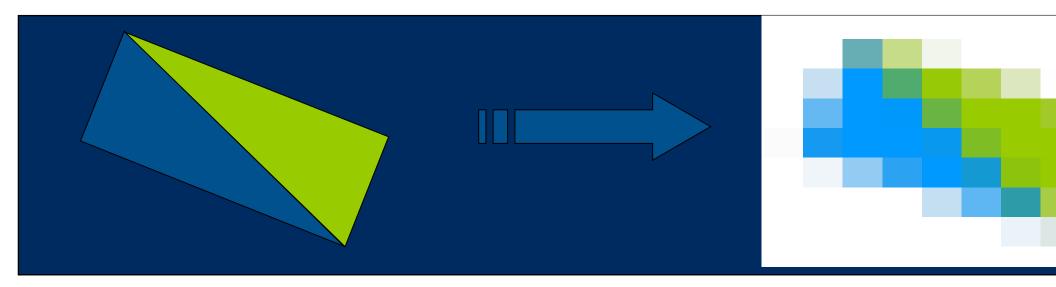  - Compute per-vertex lighting

# GPU Pipeline: Rasterizer

- ## Rasterizer
  - Convert geometric rep. (vertex) to image rep. (fragment)
    - Fragment = image fragment
      - Pixel + associated data: color, depth, stencil, etc.
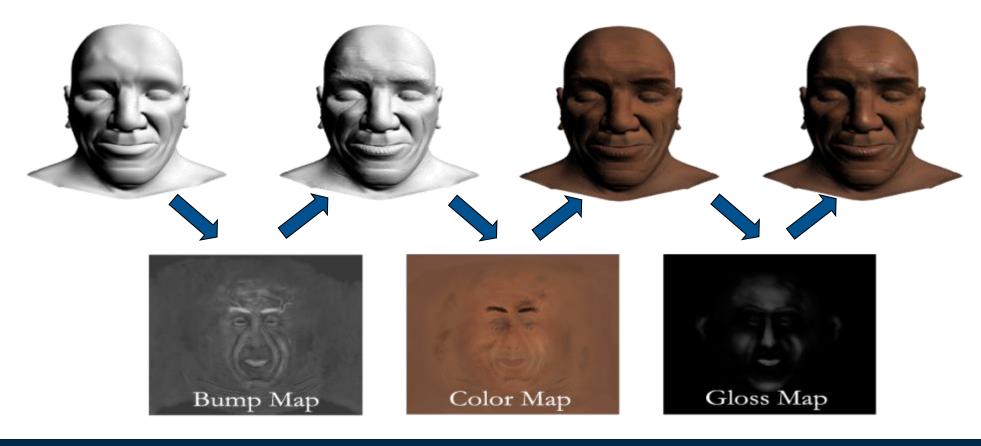  - Interpolate per-vertex quantities across pixels

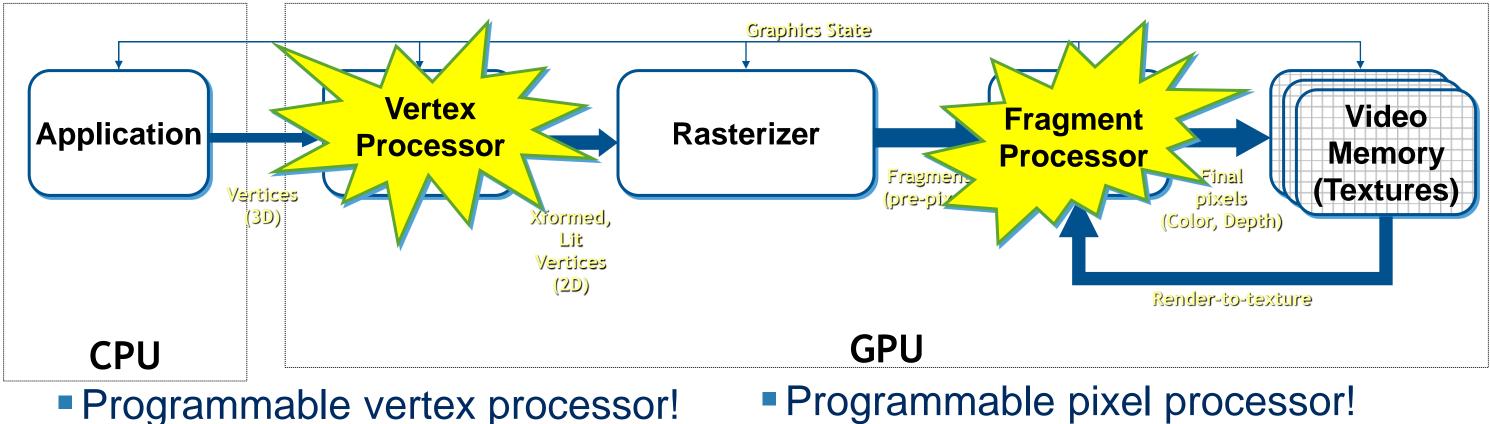- **Fragment Processors (multiple in parallel)**
  - Compute a color for each pixel
  - Optionally read colors from textures (images)



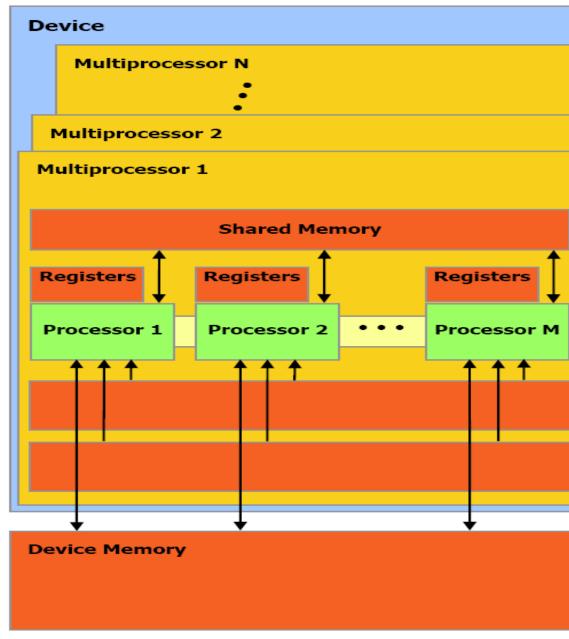Bump Map

Color Map

Gloss Map

**CPU**

**GPU**

- Programmable vertex processor!
- Programmable pixel processor!

- 16 Multiprocessors Blocks
- Each MP Block Has:
  - 8 Streaming Processors (IEEE 754 spfp compliant)
  - 16K Shared Memory
  - 64K Constant Cache
  - 8K Texture Cache
- Each processor can access all of the memory at 86Gb/s, but with different latencies:
- Shared – 2 cycle latency
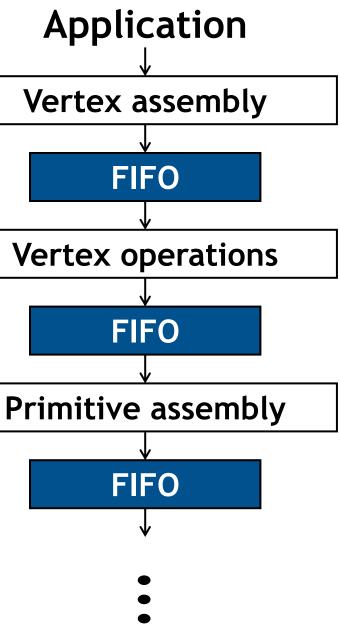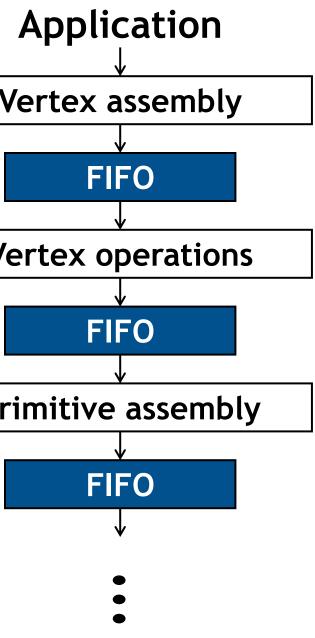- Device – 300 cycle latency

# Queueing

**FIFO buffering (first-in, first-out) is provided between task stages**
- Accommodates variation in execution time
- Provides elasticity to allow unified load balancing to work

**FIFOs can also be unified**
- Share a single large memory with multiple head-tail pairs
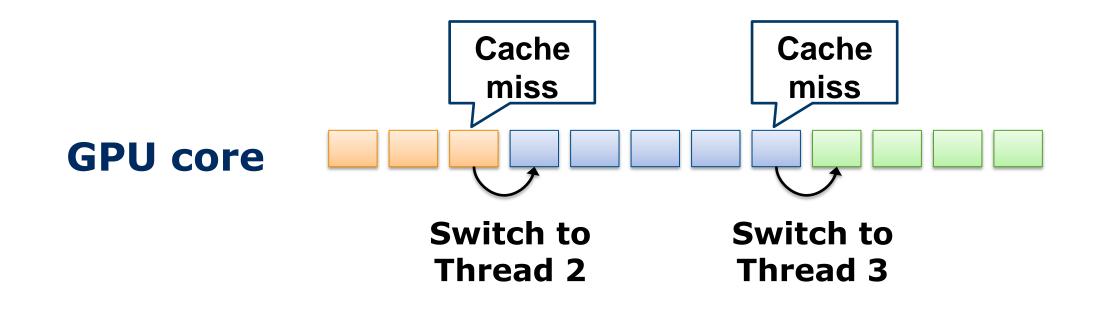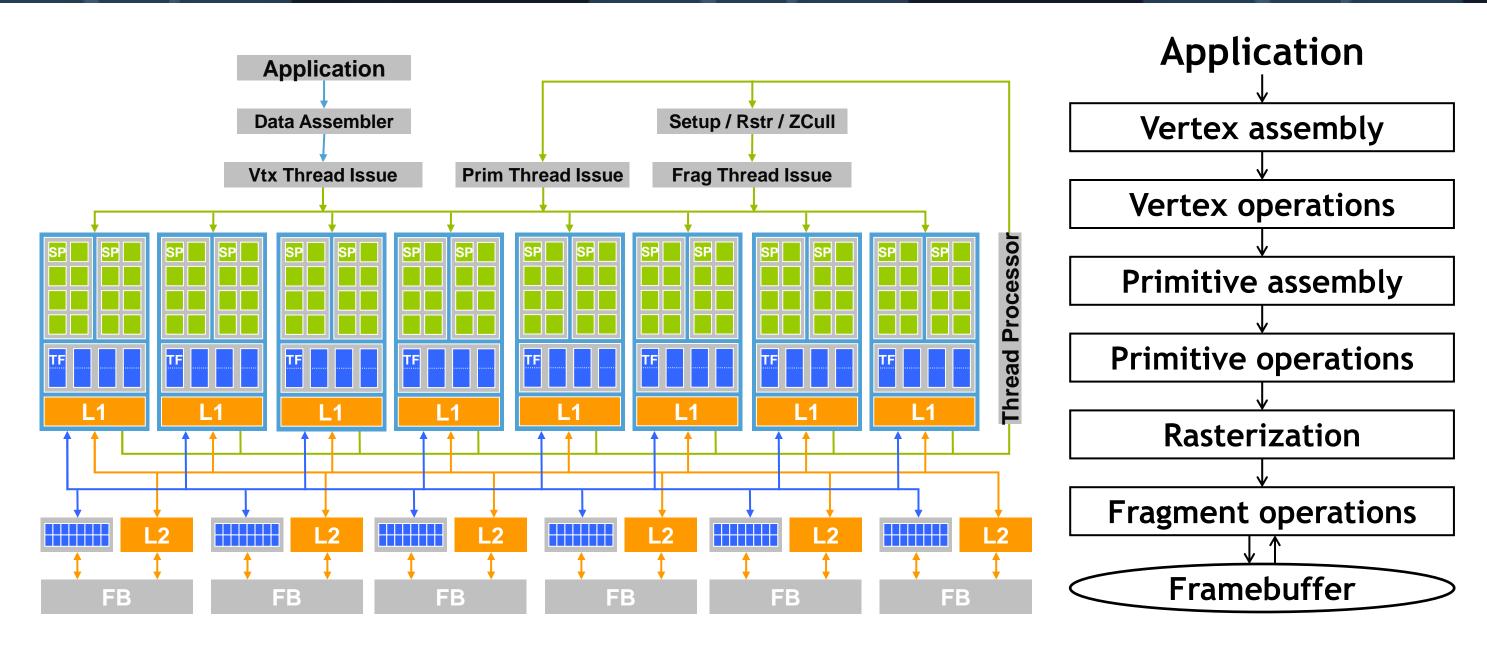- Allocate as required

**Application**

| Vertex assembly |
| FIFO |
| Vertex operations |
| FIFO |
| Primitive assembly |
| FIFO |

# GPU can effectively hide memory latency

Cache miss

Cache miss

**GPU core**

Switch to Thread 2

Switch to Thread 3

## NVIDIA GeForce 8800

## OpenGL Pipeline

Source : NVIDIA

# Correspondence (by color)

Fixed-function assembly processors

Application-programmable parallel processor

Application

| Data Assembler | *this was missing* | Setup / Rstr / ZCull |

Vtx Thread Issue    Prim Thread Issue    Frag Thread Issue

SP SP SP SP SP SP SP SP SP SP SP SP SP SP SP SP

TF  TF  TF  TF  TF  TF  TF  TF

L1  L1  L1  L1  L1

Thread Processor

Fixed-function framebuffer operations

L2  L2  L2  L2  L2  L2

FB  FB  FB  FB  FB  FB

## NVIDIA GeForce 8800

**Application**

**Vertex assembly**

**Vertex operations**

**Primitive assembly**

**Primitive operations**

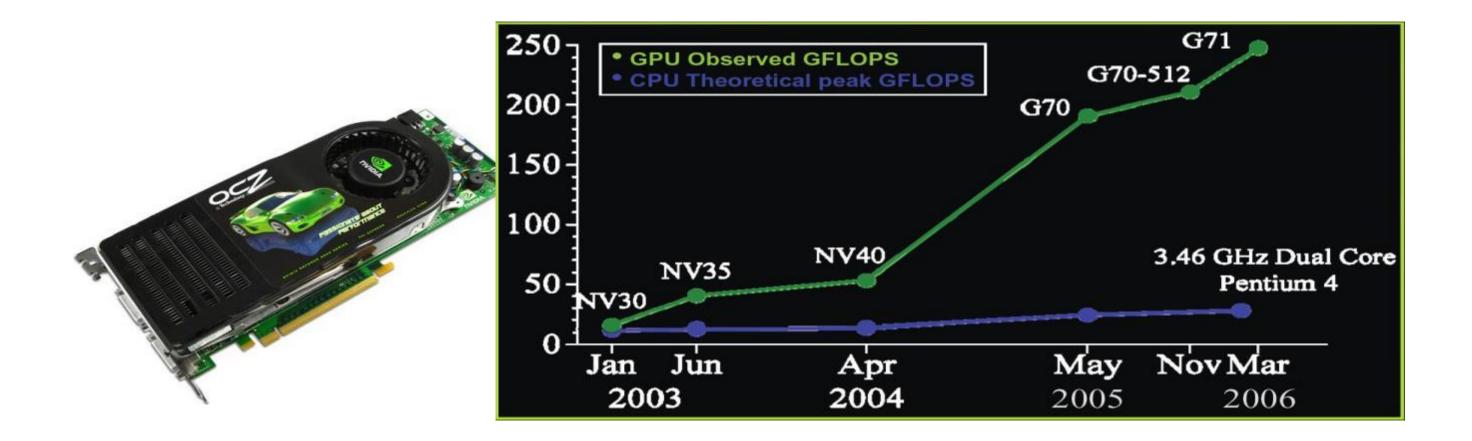**Rasterization**
(fragment assembly)

**Fragment operations**

**Framebuffer**

## OpenGL Pipeline

# The nVidia G80 GPU

- 128 streaming floating point processors @1.5Ghz
- 1.5 Gb Shared RAM with 86Gb/s bandwidth
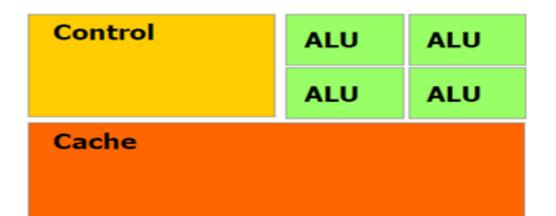- 500 Gflop on one chip (single precision)

# Why are GPU's so fast?

- Entertainment Industry has driven the economy of these chips?
  - Males age 15-35 buy $10B in video games / year
- Moore's Law ++
- Simplified design (stream processing)
  - Huge parallelism – maps well to hardware
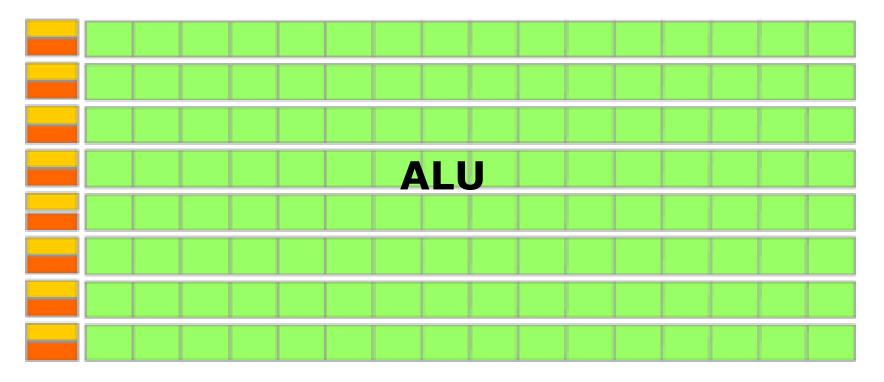  - Latency hiding using the parallelism
- Single-chip designs.

**Xeon X5550:**

**4** cores

**731M** transistors

**GTX480:**

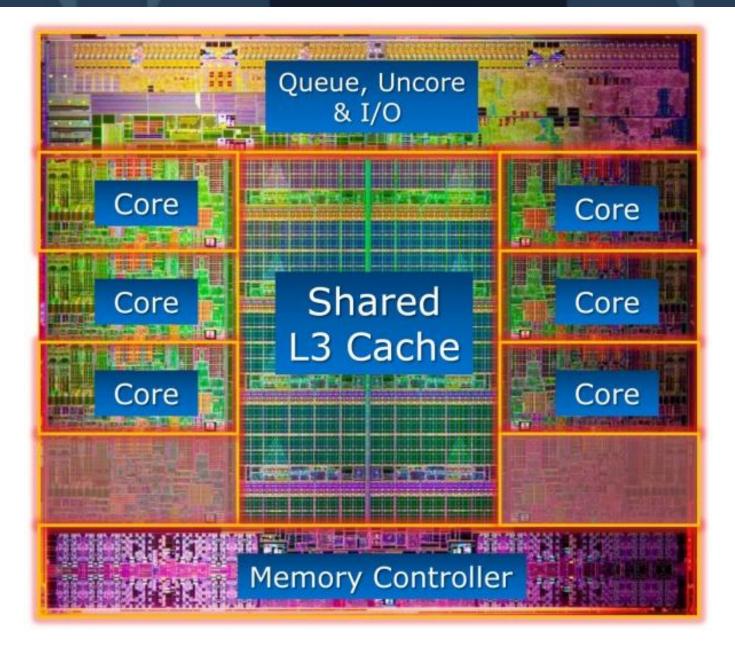**480** cores

**3,200M** transistors

**19**

CPU - Core i7

GPU – nVidia Kepler

# Very Efficient For

- Fast Parallel Floating Point Processing
- Single Instruction Multiple Data Operations
- High Computation per Memory Access

# Not As Efficient For

- Double Precision – situation is improving
- Logical Operations on Integer Data
- Branching-Intensive Operations
- Random Access, Memory-Intensive Operations

- **Programable stream processor**
  - Huge number of ALUs
  - Huge memory bandwidth
- **Programming was painful**
  - OpenGL-SL – **Shader Language**
  - Requires deep understanding of computers graphics
  - Huge applications speedup when done correctly
- **CUDA/OpenCL**
  - C-like code
  - Massively multi-threaded
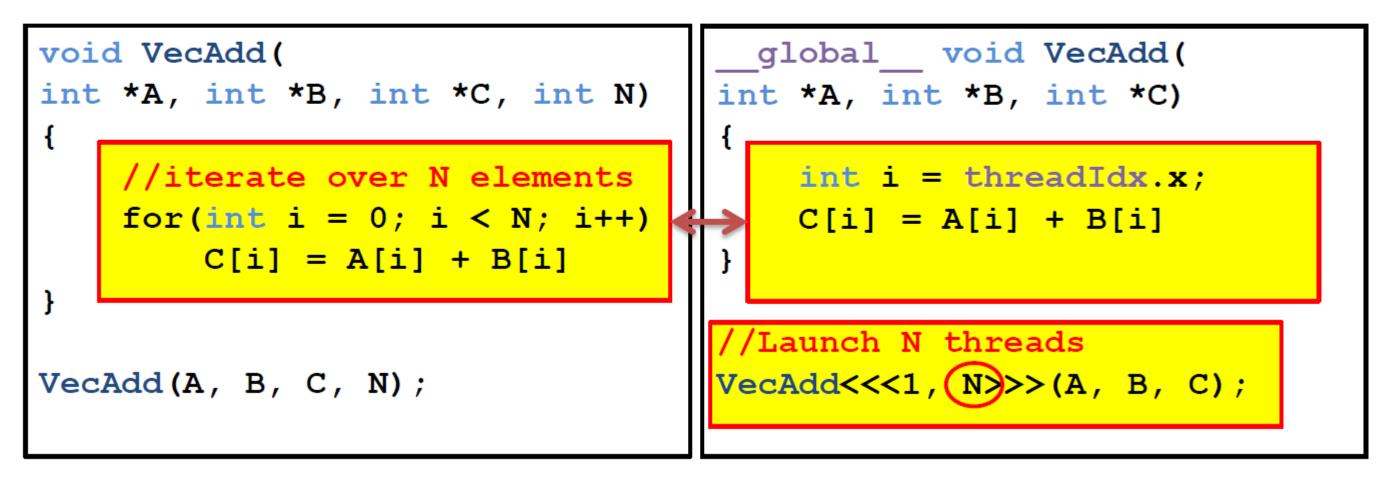  - Simple to port existing code (but not to get good performance)

## Example code: vector addition (C = A + B)

CPU code

GPU code

```
void VecAdd(
int *A, int *B, int *C, int N)
{
    //iterate over N elements
    for(int i = 0; i < N; i++)
        C[i] = A[i] + B[i]

}

VecAdd(A, B, C, N);
```

```
__global__ void VecAdd(
int *A, int *B, int *C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i]
}

//Launch N threads
VecAdd<<<1, N>>>(A, B, C);
```

# Achieving Performance in CUDA

- **Almost all C code will compile to be CUDA code**
  - But will run slower
  - Single threaded operation - ~50x slower than CPU code

- **Must expose parallelism**

- **Careful with memory accesses**
  - Thread scheduling helps hide memory access latency
  - But even this runs out

- **Moving target**
  - Performance optimizations are strongly HW and SW platform dependent

- **Can make huge difference**
  - 100x and even more

# PacketShader

## A GPU-Accelerated Software Router

# High Performance Software Router

- **Work by Sangjin Han, Keon Jang, KyoungSoo Park and Sue Moon**
  - Advanced Networking Lab, CS, KAIST
  - Networked and Distributed Computing Systems Lab, EE, KAIST
- **40 Gbps packet forwarding in a single box**
  - IPv4, 64B packets
  - Bigger packet sizes – bounded by PCI-e bandwidth
- **20 Gbps IPsec tunneling**
  - For 1024B packets
  - 10 Gbps for 64B packets

- **Despite its name, <span style="color:red">not limited to IP routing</span>**
  - You can implement whatever you want on it.

- **Driven by software**
  - Flexible
  - Friendly development environments

- **Based on commodity hardware**
  - Cheap
  - Fast evolution

# Low performance

- Due to CPU bottleneck

| Year | Ref. | H/W | IPv4 Throughput |
|------|------|-----|-----------------|
| 2008 | Egi et al. | Two quad-core CPUs | 3.5 Gbps |
| 2008 | "Enhanced SR" Bolla et al. | Two quad-core CPUs | 4.2 Gbps |
| 2009 | "RouteBricks" Dobrescu et al. | Two quad-core CPUs (2.8 GHz) | 8.7 Gbps |

- **Not capable of supporting even a single 10G port**

**Cycles needed**

IPv4 | 1,200 | + | 600 | = **1,800 cycles**
Packet I/O — IPv4 lookup

IPv6 | 1,200 | + | 1,600 | = **2,800**
Packet I/O — IPv6 lookup

IPsec | 1,200 | + | 5,400 ... | = **6,600**
Packet I/O — Encryption and hashing

**Your budget**

**1,400 cycles**
**10G, min-sized packets, dual quad-core 2.66GHz CPUs**

**(in x86, cycle numbers are from RouteBricks [Dobrescu09] and PacketShader)**
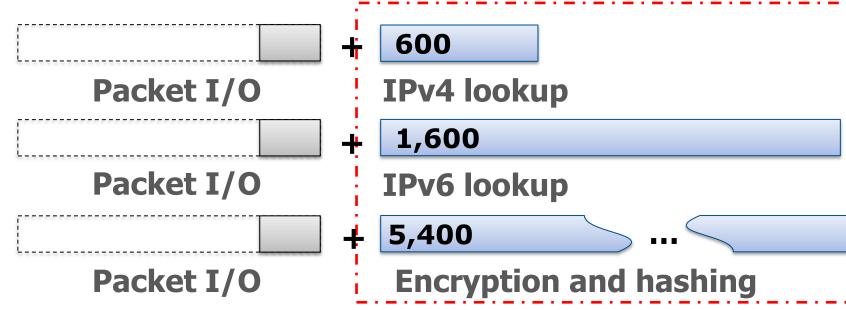
# PacketShader Approach 1: I/O Optimization

| 1,200 | + | 600 | = 1,800 cycles |
|:--:|:--:|:--:|:--|
| **Packet I/O** | | **IPv4 lookup** | |
| 1,200 | + | 1,600 | = 2,800 |
| **Packet I/O** | | **IPv6 lookup** | |
| 1,200 | + | 5,400 ... | = 6,600 |
| **Packet I/O** | | **Encryption and hashing** | |

**Packet I/O**

- **1,200 reduced to 200 cycles per packet**
- **Main ideas**
  - Huge packet buffer
  - Batch processing
- **Allocating SKBs – 50% of CPU time**

**Packet I/O** + **600**

IPv4 lookup

**Packet I/O** + **1,600**

IPv6 lookup

**Packet I/O** + **5,400** ...

**Encryption and hashing**

- **GPU Offloading for**
  - Memory-intensive or
  - Compute-intensive operations
- **Main topic of this talk**

# GPU FOR PACKET PROCESSING

1. Raw computation power

2. Memory access latency

3. Memory bandwidth

- Comparison between
  - Intel X5550 CPU
  - NVIDIA GTX480 GPU

- Compute-intensive operations in software routers
  - Hashing, encryption, pattern matching, network coding, compression, etc.
  - GPU can help!

**Instructions/sec**

**CPU: $43 \times 10^9$**
= 2.66 (GHz) ×
4 (# of cores) ×
4 (4-way superscalar)

**<**

**GPU: $672 \times 10^9$**
= 1.4 (GHz) ×
480 (# of cores)

- Software router → lots of cache misses
  - GPU can effectively hide memory latency

**CPU's memory bandwidth (theoretical): 32 GB/s**

**2. RX:**
**RAM → CPU**

**3. TX:**
**CPU → RAM**

**4. TX:**
**RAM → NIC**

**1. RX:**
**NIC → RAM**

# CPU's memory bandwidth (empirical) < 25 GB/s

**Your budget for packet processing can be less 10 GB/s**

~~Your budget for packet processing can be less 10 GB/s~~
**GPU's memory bandwidth: 174GB/s**

# HOW TO USE GPU

**Offload core operations to GPU
(e.g., forwarding table lookup)**

- For GPU, more parallelism, more throughput



**GTX480: 480 cores**

- ## The key insight
  - Stateless packet processing = parallelizable



**RX queue**

**1. Batching**

**2. Parallel Processing in GPU**

- Fast link = enough # of packets in a small time window

- 10 GbE link
  - up to 1,000 packets only in 67µs
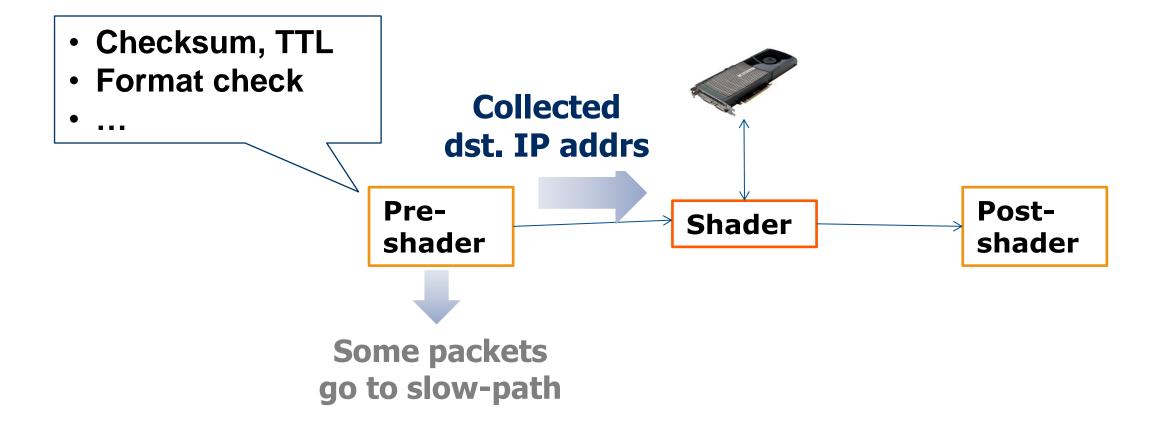
- Much less time with 40 or 100 GbE

# PACKETSHADER DESIGN

- **Three stages in a streamline**



| Pre-shader | → | Shader | → | Post-shader |

- ## IPv4 forwarding example



- **Checksum, TTL**
- **Format check**
- **…**

**Collected dst. IP addrs**

**Pre-shader**

**Shader**

**Post-shader**

**Some packets go to slow-path**

- IPv4 forwarding example

**2. Forwarding table lookup**

**1. IP addresses**  **3. Next hops**

| Pre-shader | Shader | Post-shader |

- IPv4 forwarding example

# Scaling with a Multi-Core CPU



**Master core**

| Device driver | Pre-shader | Shader | Post-shader | Device driver |

**Worker cores**

# EVALUATION

**CPU:**



Total 8 CPU cores

**Quad-core, 2.66 GHz**

**NIC:**



Total 80 Gbps

**Dual-port 10 GbE**

**GPU:**



Total 960 cores

**480 cores, 1.4 GHz**

Input traffic

Processed packets

**8 × 10 GbE links**

# Packet generator
**(Up to 80 Gbps)**

# PacketShader

# Results (w/ 64B packets)
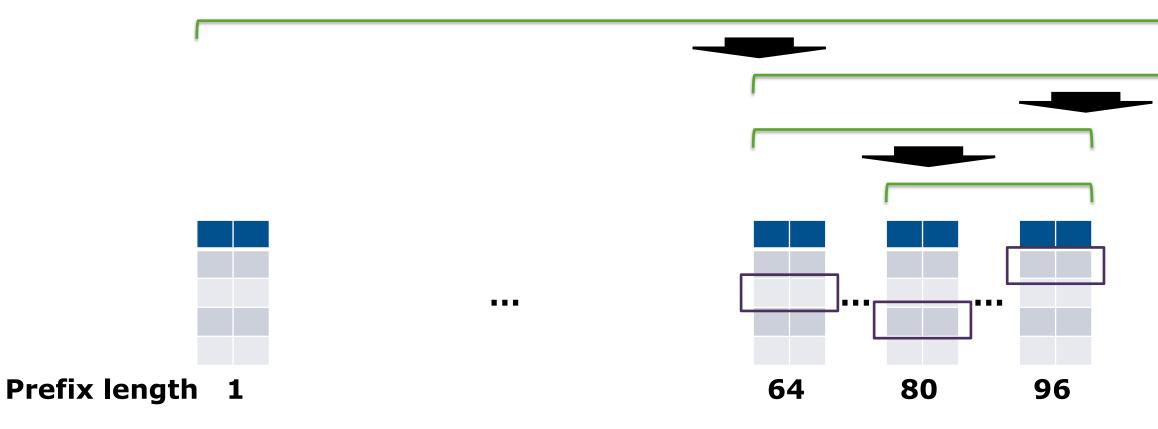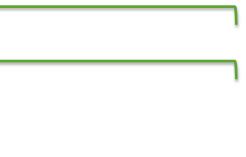
- Longest prefix matching on 128-bit IPv6 addresses

- Algorithm: binary search on hash tables [Waldvogel97]
  - 7 hashings + 7 memory accesses


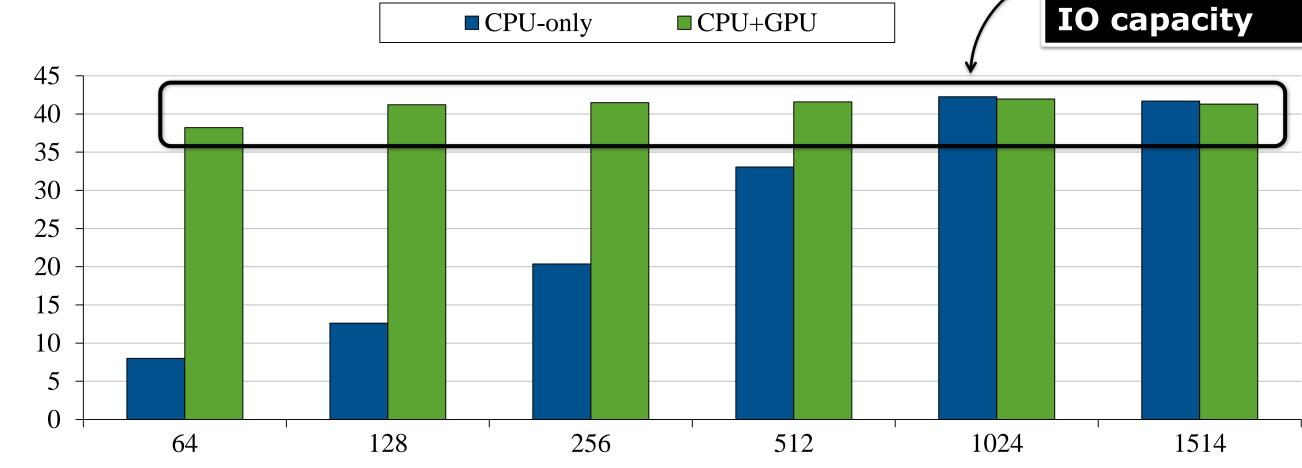
**Prefix length**    1                                        64        80        96                128

# Example 1: IPv6 forwarding



**Bounded by motherboard IO capacity**

Legend: ■ CPU-only  ■ CPU+GPU

X-axis: Packet size (bytes) — 64, 128, 256, 512, 1024, 1514

Y-axis: Throughput (Gbps)

(Routing table was randomly generated with 200K entries)
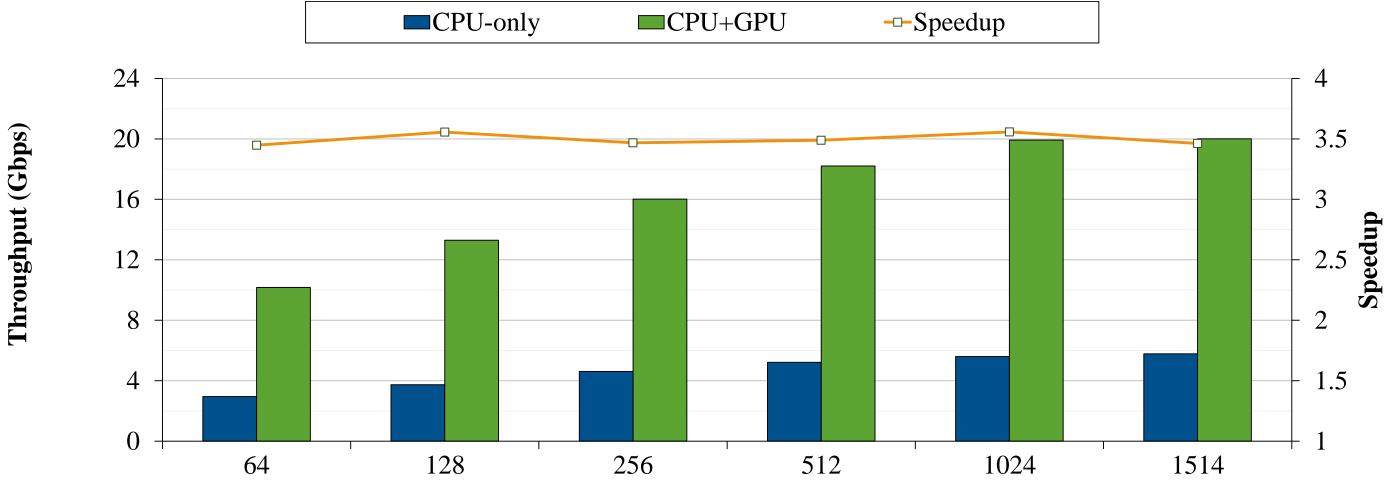
- **ESP (Encapsulating Security Payload) Tunnel mode**
  - with AES-CTR (encryption) and SHA1 (authentication)

- 3.5x speedup

| Year | Ref. | H/W | IPv4 Throughput | |
|------|------|-----|----------------:|---|
| 2008 | Egi *et al.* | Two quad-core CPUs | 3.5 Gbps | |
| 2008 | "Enhanced SR" Bolla *et al.* | Two quad-core CPUs | 4.2 Gbps | **Kernel** |
| 2009 | "RouteBricks" Dobrescu *et al.* | Two quad-core CPUs (2.8 GHz) | 8.7 Gbps | |
| 2010 | PacketShader (CPU-only) | Two quad-core CPUs (2.66 GHz) | 28.2 Gbps | **User** |
| 2010 | PacketShader (CPU+GPU) | Two quad-core CPUs + two GPUs | 39.2 Gbps | |

# GPU

- a great opportunity for fast packet processing

# PacketShader

- Optimized packet I/O + GPU acceleration
- scalable with
  - \# of multi-core CPUs, GPUs, and high-speed NICs

# Current Prototype

- Supports IPv4, IPv6, OpenFlow, and IPsec
- 40 Gbps performance on a single PC

- **Control plane integration**
  - Dynamic routing protocols with Quagga or Xorp

- **Multi-functional, modular programming environment**
  - Integration with Click? [Kohler99]

- **Opportunistic offloading**
  - CPU at low load
  - GPU at high load

- **Stateful packet processing**

# SSLShader

## A GPU-Accelerated Software Router

Thank You

Mellanox® TECHNOLOGIES

Connect. Accelerate. Outperform.™