**Parul**® University
Vadodara, Gujarat

NAAC A++
GRADE

Information and
Communication Technology

# 5:Parallel Programming Using GPU

**Kusum Lata Dhiman**

**Assistant Professor**

Computer Science & Engineering

## Module Topics

**Parallel Programming Using GPGPU:** An Overview of GPGPU, DGX architecture, An Overview of GPGPU Programming, An Overview of GPGPU Memory Hierarchy Features, CUDA Programming

# Introduction to GPGPU

- GPGPU, or General-Purpose computing on Graphics Processing Units, is a technology that leverages the immense parallel processing power of graphics processing units (GPUs) for tasks beyond traditional graphics rendering. GPUs were initially designed to handle the complex calculations required for rendering images and video in computer graphics, but they are well-suited for a wide range of general-purpose computing tasks due to their highly parallel architecture.

# Introduction to GPGPU

- **1. Parallel Processing Power:** GPUs consist of thousands of small processing cores that can perform multiple calculations simultaneously. This parallel architecture makes GPUs exceptionally well-suited for tasks that involve performing the same operation on a large dataset or solving complex mathematical problems.

- **2. Transition from Graphics to General-Purpose:** The transition from using GPUs solely for graphics-related tasks to general-purpose computing began in the early 2000s. Researchers and developers realized that GPUs could be used for a wide variety of computational tasks, including scientific simulations, data analytics, machine learning, and more.

# Introduction to GPGPU

- **3. Programming Models:** To harness the power of GPUs for general-purpose computing, specialized programming models and libraries were developed. CUDA (Compute Unified Device Architecture) by NVIDIA and OpenCL (Open Computing Language) are two common frameworks that allow developers to write code for GPUs.

- **4. Parallelism:** GPGPU programming often involves breaking down a problem into smaller parallelizable tasks that can be executed simultaneously on the GPU cores. This approach can significantly accelerate computations compared to running the same code on a traditional CPU.

# Introduction to GPGPU

- **5. Applications:** GPGPU is widely used in various fields, including:
- **Scientific Computing:** GPGPU accelerates simulations in physics, chemistry, and engineering.
- **Data Analytics:** GPUs speed up data processing, especially for tasks like machine learning and data mining.
- **Computer Vision:** GPUs enhance image and video analysis.
- **Cryptocurrency Mining:** GPUs are used to mine cryptocurrencies due to their computational power.
- **Gaming:** Besides their primary role in rendering graphics, GPUs are also used for physics simulations and AI-driven game mechanics.

# Introduction to GPGPU

**6. Challenges:** While GPGPU computing offers significant performance advantages, it also presents challenges. Developers need to learn specialized programming languages, manage data transfer between the CPU and GPU, and optimize algorithms for parallel execution.

**7. Future Developments:** GPGPU technology continues to evolve. Newer GPUs are designed with features specifically for AI and machine learning tasks, and hardware improvements, such as ray tracing capabilities, are also becoming important in graphics and simulation applications.

**Parul**® University
Vadodara, Gujarat

NAAC
GRADE A++

Information and
Communication Technology

# DGX architecture

**Nvidia DGX** is a line of Nvidia-produced servers and workstations which specialize in using GPGPU to accelerate deep learning applications. The typical design of a DGX system is based upon a rackmount chassis with motherboard that carries high performance x86 server CPUs (Typically Intel Xeons, with the exception DGX A100 and DGX Station A100, which both utilize AMD EPYC CPUs). The main component of a DGX system is a set of 4 to 16 Nvidia Tesla GPU modules on an independent system board. DGX systems have large heatsinks and powerful fans to adequately cool thousands of watts of thermal output. The GPU modules are typically integrated into the system using a version of the SXM socket.

**Parul**® University
Vadodara, Gujarat

NAAC
GRADE A++

Information and
Communication Technology

# Key components and features of DGX

- **Multiple GPUs:** DGX systems are known for their inclusion of multiple high-end NVIDIA GPUs, typically based on the latest GPU architecture available at the time. These GPUs are optimized for AI workloads and are tightly integrated into the system.

- **NVLink Interconnect:** NVLink is NVIDIA's high-speed GPU interconnect technology that allows for faster data transfer and communication between GPUs within the same system. It enables better scaling for multi-GPU deep learning tasks.

**Parul**®University
Vadodara, Gujarat

NAAC
GRADE A++

Information and
Communication Technology

## Key components and features of DGX

- **High-speed Networking:** DGX systems often feature high-speed networking capabilities, such as InfiniBand or Ethernet, to enable fast data transfer and communication between multiple DGX nodes in a cluster. This is crucial for distributed AI training.

- **Optimized Cooling and Form Factor:** DGX systems are designed to provide efficient cooling for the powerful GPUs and other components, ensuring consistent and reliable performance. They often come in rack-mounted form factors suitable for data center deployment.

**Parul**® University
Vadodara, Gujarat

NAAC
GRADE A++

Information and
Communication Technology

## Key components and features of DGX

- **AI Software Stack:** DGX systems are typically bundled with a comprehensive software stack that includes NVIDIA's CUDA platform, cuDNN (NVIDIA's deep learning library), and other AI-specific libraries and frameworks. This software stack is pre-configured and optimized for AI workloads.

- **Storage Solutions:** DGX systems may offer various storage options, including high-speed SSDs and NVMe storage, to support the large datasets commonly used in deep learning tasks.

# Key components and features of DGX

- **Deep Learning Performance:** DGX systems are benchmarked and optimized for deep learning tasks, making them ideal for training large neural networks quickly and efficiently.

- **GPU Cloud Integration:** NVIDIA may offer cloud-based solutions that integrate with DGX systems, allowing users to access additional GPU resources and scale their AI workloads as needed.

# overview of GPGPU programming

**1. Parallelism:** At the heart of GPGPU programming is the concept of parallelism. GPUs consist of thousands of small processing cores that can perform calculations simultaneously. Unlike CPUs, which are optimized for sequential processing, GPUs excel at parallel processing tasks where the same operation is performed on multiple pieces of data simultaneously. This parallelism is key to achieving significant speedup for certain types of computations.

# overview of GPGPU programming

**2. Programming Models:**

**CUDA (Compute Unified Device Architecture):** Developed by NVIDIA, CUDA is one of the most widely used programming models for GPGPU programming. It provides a set of extensions to the C and C++ programming languages, allowing developers to write GPU-accelerated code. CUDA provides tools for managing data transfers between the CPU and GPU and for launching parallel kernels (functions executed on the GPU).

# overview of GPGPU programming

**OpenCL (Open Computing Language):** OpenCL is an open standard for GPGPU programming supported by various hardware vendors, including NVIDIA, AMD, and Intel. It allows developers to write code that can run on different GPU architectures. OpenCL provides a lower-level programming model compared to CUDA and is suitable for more platform-independent GPGPU applications.

# overview of GPGPU programming

**3. Data Transfer:** GPGPU programming often involves transferring data between the CPU and GPU memory. Efficient management of data transfers is crucial for achieving good performance. Techniques like pinned memory, asynchronous data transfers, and memory coalescing are used to optimize data movement.

**4. Kernel Execution:** In GPGPU programming, code that runs on the GPU is referred to as a "kernel." Kernels are designed to be highly parallelizable and are executed across multiple GPU threads or cores. Developers need to carefully design and optimize their kernels to maximize GPU utilization.

# overview of GPGPU programming

**5. Thread Synchronization:** Handling synchronization and coordination between threads running on the GPU is essential in GPGPU programming. Developers must be aware of techniques like thread synchronization barriers and atomic operations to ensure correctness and avoid race conditions.

# overview of GPGPU programming

**6. Applications:**

**Scientific Computing:** GPGPU is extensively used for scientific simulations, including physics, chemistry, and climate modeling.

**Machine Learning and Deep Learning:** Training and inference of neural networks are accelerated using GPUs, contributing to the success of AI applications.

**Data Analysis:** GPUs accelerate data processing tasks like data transformation, filtering, and statistical analysis.

**Image and Video Processing:** GPGPU is used in image and video manipulation, including image filtering, compression, and computer vision tasks.

# overview of GPGPU programming

**7. Challenges:** GPGPU programming presents challenges such as a steeper learning curve, managing GPU resources, debugging and profiling GPU code, and ensuring data consistency.

**8. Performance Optimization:** Achieving optimal performance in GPGPU programming often requires optimizing memory access patterns, minimizing data transfers between the CPU and GPU, and fine-tuning kernel code.

**Parul**® University
Vadodara, Gujarat

NAAC
GRADE A++

Information and
Communication Technology

# Overview of GPGPU Memory Hierarchy Features

GPUs have a memory hierarchy that consists of different types of memory with varying characteristics and access speeds. Here's an overview of the key memory hierarchy features in GPGPU:

**1. Global Memory:**

- **Characteristics:** Global memory is the largest and slowest memory in the GPU's memory hierarchy. It is typically implemented as off-chip DRAM (Dynamic Random-Access Memory).

- **Usage:** Global memory is used for storing data that needs to be shared among all threads in a GPU kernel. It is the primary storage for input data, intermediate results, and output data.

# Overview of GPGPU Memory Hierarchy Features

- **Access:** Access to global memory has higher latency and is typically slower than other memory types. It requires careful management to minimize memory access times.

## 2. Shared Memory:

- **Characteristics:** Shared memory is a small, fast, on-chip memory that is shared among threads within a thread block (also known as a warp). It provides low-latency access and is used for thread communication and for creating thread-level software-managed caches.

- **Usage:** Shared memory is used for storing data that is frequently accessed by threads within the same thread block. It can significantly reduce memory access times when used effectively.

# Overview of GPGPU Memory Hierarchy Features

- **Access:** Access to shared memory is much faster than global memory but is limited by its size, which is typically small compared to global memory.

## 3. Local Memory:

- **Characteristics:** Local memory is similar to global memory but is specific to each thread. It is often implemented in global memory, but it is used when a thread's data does not fit in registers or shared memory.

- **Usage:** Local memory is used for thread-specific data that cannot be stored in registers or shared memory. It is usually accessed when registers are exhausted.

## Overview of GPGPU Memory Hierarchy Features

- **Access:** Access to local memory is slower than registers or shared memory but faster than global memory.

**4. Registers:**

- **Characteristics:** Registers are the fastest and most abundant on-chip memory in the GPU. Each thread has its set of registers for storing variables and temporary data.

- **Usage:** Registers are used for storing thread-specific variables, loop counters, and temporary data needed for computations.

- **Access:** Access to registers is extremely fast, and they are used for minimizing memory latency.

## Overview of GPGPU Memory Hierarchy Features

**5. Texture Memory and Constant Memory:**

- **Characteristics:** Texture memory and constant memory are read-only memory types optimized for specific access patterns. Texture memory is designed for 2D spatial locality, while constant memory is for read-only constant data.

- **Usage:** These memory types are used when read-only data with specific access patterns can be exploited to improve memory access efficiency.

- **Access:** Access to texture memory and constant memory is optimized for specific access patterns and can be faster than global memory.

## Overview of GPGPU Memory Hierarchy Features

## 6. L1 and L2 Cache:

- **Characteristics:** Many modern GPUs have L1 and L2 caches that sit between the on-chip memory hierarchy and global memory. These caches help reduce the latency of memory accesses by storing frequently accessed data.

- **Usage:** Caches are used to automatically manage data movement between the memory hierarchy levels and global memory.

- **Access:** Access to data in caches is faster than accessing global memory directly.

# Overview of CUDA Programming

## 1. CUDA Programming Model:

- **Host and Device:** In CUDA programming, there are two primary components: the host, which represents the CPU and runs the main program, and the device, which represents the GPU and executes parallel kernels.

- **Parallelism:** The core concept of CUDA is parallelism. Developers write kernels, which are functions that can be executed in parallel by multiple threads on the GPU. These threads work together to perform computations on data.

Parul® University
Vadodara, Gujarat

NAAC
GRADE A++

Information and
Communication Technology

## Overview of CUDA Programming

**Hierarchy:** CUDA programs are organized into grids, blocks, and threads. A grid is composed of multiple blocks, and each block consists of multiple threads. This hierarchy allows for fine-grained control of parallel execution.

**2. CUDA Language Extensions:**

- **CUDA C/C++:** CUDA extends the C/C++ programming languages with special keywords and constructs for specifying parallelism and managing data transfer between the CPU and GPU.

- **Kernel Functions:** Kernels are functions that are executed in parallel by many threads on the GPU. They are marked with the **__global__** keyword in CUDA C/C++

## Overview of CUDA Programming

- **Thread IDs:** CUDA provides built-in variables (**threadIdx**, **blockIdx**, and **blockDim**) that allow threads to identify their position within the grid and block structure.

- **Synchronization:** CUDA provides synchronization mechanisms such as barriers (**__syncthreads()**) to coordinate the execution of threads within a block.

## Overview of CUDA Programming

**3. Memory Management:**

- **Global Memory:** Global memory is the primary storage for data that is accessible by both the host and the device. Efficient memory management, including data transfers between global memory and CPU memory, is crucial for performance.

- **Shared Memory:** Shared memory is a fast, on-chip memory that allows threads within the same block to share data. It is used for communication and as a programmer-managed cache.

- **Constant and Texture Memory:** These memory types are optimized for specific access patterns and can be used to improve memory access efficiency.

**4. Data Transfer**

•**Memcpy:** Data transfer between CPU and GPU memory is often done using CUDA-aware memory copy functions (**cudaMemcpy**).

•**Unified Memory:** CUDA offers Unified Memory, which simplifies memory management by allowing a single memory space that is accessible by both the CPU and GPU. The system manages data migration between CPU and GPU memory transparently.

**Parul**® University
Vadodara, Gujarat

NAAC A++
GRADE

Information and
Communication Technology

## Overview of CUDA Programming

**5. Performance Optimization:**

- **Thread Divergence:** Minimizing thread divergence (variations in execution paths) is essential for efficient parallel execution.

- **Memory Access Patterns:** Optimizing memory access patterns and minimizing global memory access latency are critical for performance.

- **Shared Memory Usage:** Efficient utilization of shared memory and avoiding bank conflicts are important for maximizing throughput.

**Parul**® University
Vadodara, Gujarat

NAAC **A**++
GRADE

Information and
Communication Technology

## Overview of CUDA Programming

**6. Tools and Libraries:**

- NVIDIA provides a suite of development tools, including the CUDA Toolkit, which includes compilers, debuggers, and profilers.

- Various CUDA libraries, such as cuBLAS (for linear algebra), cuDNN (for deep learning), and cuFFT (for Fast Fourier Transforms), accelerate common mathematical and scientific operations.

**Parul®**
University

**NAAC**
**GRADE** **A++**

https://paruluniversity.ac.in/