## Department of Computer Science and Engineering, Indian Institute of Technology Madras

| Course title | **GPU Programming** | | | | | | | | Course No | **CSNEW** | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Department | CSE | New Credits | L | T | E | P | O | C | Old Credits | L | T | P | C |
| | | | 3 | 1 | 0 | 0 | 8 | 12 | | 3 | 1 | 0 | 4 |
| Offered for | BTech / DD / MTech / MS / PhD | | | | | | | | Status | Final | | | |
| Faculty | **Rupesh Nasre.** | | | | | | | | Type | Theory | | | |
| Pre-requisite | | | | | | | | | To take effect from | 2017 | | | |
| Submission date | Date of approval by DCC | | | Date of approval by BAC | | | | | Date of approval by Senate | | | | |
| | | | | | | | | | | | | | |

**Objectives:**
To learn parallel programming with graphics processing units (GPUs)

**Outcomes:**
Students would learn concepts in parallel programming, implementation of programs on GPUs, debugging and profiling parallel programs.

**Course Contents:**

*Topic (number of lectures + number of tutorials)*

Introduction (2 + 1):
  - history, graphics processors, graphics processing units, GPGPUs
  - clock speeds, CPU / GPU comparisons, heterogeneity
  - accelerators, parallel programming, CUDA / OpenCL / OpenACC, Hello World

Computation (3 + 1):
  - kernels, launch parameters
  - thread hierarchy, warps / wavefronts, thread blocks / workgroups, streaming multiprocessors
  - 1D / 2D / 3D thread mapping, device properties, simple programs

Memory (8 + 2):
  - memory hierarchy, DRAM / global, local / shared, private / local, textures, constant memory
  - pointers, parameter passing, arrays and dynamic memory, multi-dimensional arrays
  - memory allocation, memory copying across devices
  - programs with matrices, performance evaluation with different memories

Synchronization (6 + 2):
  - memory consistency
  - barriers (local versus global), atomics, memory fence
  - prefix sum, reduction
  - programs for concurrent data structures such as worklists, linked-lists
  - synchronization across CPU and GPU

Functions (3 + 1):
  - device functions, host functions, kernels, functors
  - using libraries (such as Thrust), developing libraries

Support (1 + 2):
  - debugging GPU programs
  - profiling, profile tools, performance aspects

Streams (3 + 1):
  - asynchronous processing, tasks, task-dependence
  - overlapped data transfers, default stream, synchronization with streams
  - events, event-based-synchronization

- overlapping data transfer and kernel execution, pitfalls

Case studies (3 + 2):
 - image processing
 - graph algorithms
 - simulations
 - deep learning
 - ...

Advanced topics (8 + 2):
 - dynamic parallelism
 - unified virtual memory
 - multi-GPU processing
 - peer access
 - heterogeneous processing

Course evaluation would involve programming assignments.

**Text Books:**
Programming Massively Parallel Processors: A Hands-on Approach; David Kirk, Wen-mei Hwu; Morgan Kaufman; 2010 (ISBN: 978-0123814722)

**Reference Books:**

CUDA Programming: A Developer's Guide to Parallel Computing with GPUs; Shane Cook; Morgan Kaufman; 2012 (ISBN: 978-0124159334)