# Syllabus

# 410241 : High Performance Computing

| Teaching Scheme : | Credit : 04 | Examination Scheme : |
|---|---|---|
| TH : 04 Hours/Week | | In-Sem (Paper) : 30 Marks |
| | | End-Sem (Paper) : 70 Marks |

## Prerequisite Courses :

210253-Microprocessor, 210244 - Computer Organization and Architecture, 210254-Principles of Programming Languages, 310251- Systems Programming and Operating System

## Course Objectives :

- To study parallel computing hardware and programming models
- To be conversant with performance analysis and modeling of parallel programs
- To understand the options available to parallelize the programs
- To know the operating system requirements to qualify in handling the parallelization

## Course Outcomes :

On completion of the course, student will be able to

- Describe different parallel architectures, inter-connect networks, programming models
- Develop an efficient parallel algorithm to solve given problem
- Analyze and measure performance of modern parallel computing systems
- Build the logic to parallelize the programming task

## Course Contents

### Unit I : Introduction                                            (09 Hours)

Motivating Parallelism, Scope of Parallel Computing, Parallel Programming Platforms : Implicit Parallelism, Trends in Microprocessor and Architectures, Limitations of Memory, System Performance, Dichotomy of Parallel Computing Platforms, Physical Organization of Parallel Platforms, Communication Costs in Parallel Machines, Scalable design principles, Architectures : N-wide superscalar architectures, Multi-core architecture.

**(Refer Chapter 1)**

### Unit II : Parallel Programming                                    (09 Hours)

Principles of Parallel Algorithm Design : Preliminaries, Decomposition Techniques, Characteristics of Tasks and Interactions, Mapping Techniques for Load Balancing, Methods for Containing Interaction Overheads, Parallel Algorithm Models, The Age of Parallel Processing, the Rise of GPU Computing, A Brief History of GPUs, Early GPU.

**(Refer Chapter 2)**

## Unit III : Basic Communication
(09 Hours)

Operations- One-to-All Broadcast and All-to-One Reduction, All-to-All Broadcast and Reduction, All-Reduce and Prefix-Sum Operations, Scatter and Gather, All-to-All Personalized Communication, Circular Shift, Improving the Speed of Some Communication Operations.
**(Refer Chapter 3)**

## Unit IV : Analytical Models of Parallel Programs
(9 Hours)

Analytical Models : Sources of overhead in Parallel Programs, Performance Metrics for Parallel Systems, and The effect of Granularity on Performance, Scalability of Parallel Systems, Minimum execution time and minimum cost, optimal execution time. Dense Matrix Algorithms : Matrix-Vector Multiplication, Matrix-Matrix Multiplication.
**(Refer Chapter 4)**

## Unit V : Parallel Algorithms- Sorting and Graph
(09 Hours)

Issues in Sorting on Parallel Computers, Bubble Sort and its Variants, Parallelizing Quick sort, All-Pairs Shortest Paths, Algorithm for sparse graph, Parallel Depth-First Search, Parallel Best-First Search. **(Refer Chapter 5)**

## Unit VI : CUDA Architecture
(09 Hours)

CUDA Architecture, Using the CUDA Architecture, Applications of CUDA Introduction to CUDA C-Write and launch CUDA C kernels, Manage GPU memory, Manage communication and synchronization, Parallel programming in CUDA- C.
**(Refer Chapter 6)**

❑❑❑

## UNIT I

### Chapter 1 : Parallel Processing Concepts    1-1 to 1-44

**Syllabus :** Motivating Parallelism, Scope of Parallel Computing, Parallel Programming Platforms : Implicit Parallelism, Trends in Microprocessor and Architectures, Limitations of Memory, System Performance, Dichotomy of Parallel Computing Platforms, Physical Organization of Parallel Platforms, Communication Costs in Parallel Machines, Scalable design principles, Architectures : N-wide superscalar architectures, Multi-core architecture.

### UNIT II

**Chapter 2 :    Parallel Programming        2-1 to 2-25**

> **Syllabus :** Principles of Parallel Algorithm Design : Preliminaries, Decomposition Techniques, Characteristics of Tasks and Interactions, Mapping Techniques for Load Balancing, Methods for Containing Interaction Overheads, Parallel Algorithm Models, The Age of Parallel Processing, the Rise of GPU Computing, A Brief History of GPUs, Early GPU.

## UNIT V

## Chapter 5 :   Parallel Algorithms- Sorting and Graph
### 5-1 to 5-21

**Syllabus :** Issues in Sorting on Parallel Computers, Bubble Sort and its Variants, Parallelizing Quick sort, All-Pairs Shortest Paths, Algorithm for sparse graph, Parallel Depth-First Search, Parallel Best-First Search.

## UNIT VI

## Chapter 6 :   CUDA Architecture     6-1 to 6-23

**Syllabus :** CUDA Architecture, Using the CUDA Architecture, Applications of CUDA Introduction to CUDA C-Write and launch CUDA C kernels, Manage GPU memory, Manage communication and synchronization, Parallel programming in CUDA-C.

❑❑❑