

Lectures

Introduction

L0 and L1

Program Parallelization

Decomposition: Decompose a sequential algorithm into tasks (programmer)

- Granularity of tasks are important
- Tasks have dependencies (data or control) between each other which defines the execution order

Scheduling: Assign tasks to processes (programmer / compiler)

Mapping - Map processes to cores (OS)

Von Neumann Computation Model instruction and data are stored in memory, and processors computes.

Memory Wall disparity between memory speed and processor speed ($\leq 1 \text{ ns}$ VS $\geq 100 \text{ ns}$)

Processing unit refers to a core that can execute a kernel thread

Interconnect busses between different components in the machine

Node Machine in a distributed system

Why Parallel

Primary Reasons

- 1 OVercome limits of serial computing
- 2 Solve larger problems
- 3 Save (wall-clock) time

Other Reasons

- Take advantage of non-local resources
- Cost/energy saving - use multiple cheaper computing resources
- Overcome memory constraints

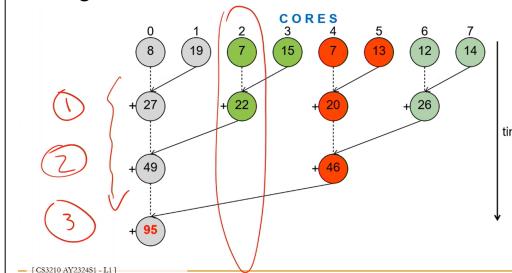
Computational Model Attributes

- **Operation mechanism** Primitive units of computation or basic actions of the computer on a specific Architecture
- **Data Mechanism** How we access and store data in address space
- **Control Mechanism** How primitive units of computation are scheduled
- **Communication Mechanism** Modes and patterns of exchanging information between parallel tasks (e.g message passing, shared memory)
- **Synchronization Mechanism** ensures to ensure needed information arrives at the right time

Dependencies and Coordination

- Dependencies among tasks impose constraints on scheduling
- Memory organizations: Shared-memory (threads), distributed-memory (processes)
- Coordination (synchronization) imposes additional overheads

Two algorithms



- Core 0 is active throughout the execution
- Some cores are idle
- This is a lot better than having all cores idle while the master core is executing

Parallel Performance

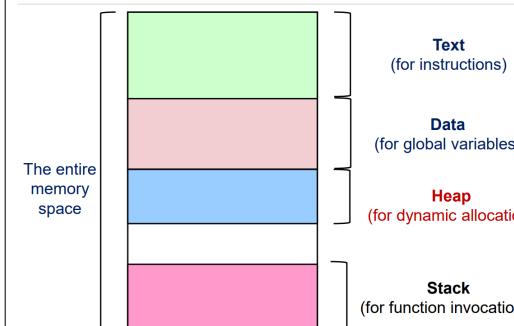
- Execution time Vs Throughput
- Parallel execution time = computation time + parallelization overheads
- Overheads: Distribution of work(tasks) to processes, information exchange, synchronisation, idle time, etc

Background on Parallelism

L2: Processes and Threads

Process

- Identified by PID
- Program counter, global data (open files, network connections), stack or heap, current values of the registers (GPRs and Special)
- These information are abstracted in the PCB, and each process can be viewed as having exclusive access to its address space
- Explicit communication is needed
- **Disadvantage**
 1. High overhead of system calls
 2. Potential re-allocation of data-structures
 3. Communication goes through OS (system calls) and context switch is costly



Multi tasking

- Overhead: Context switching (PCB change) is needed and states of suspended process must be saved
- Time slicing: Pseudo-parallelism
- Child processes can use parent's data
- **Inter-process communication (IPC)**
- Shared memory: need to protect access with locks

- Message passing: Blocking, unblocking, Synchronous, unsynchronous

Exceptions

- Executing a **machine level instruction** can cause exception
- For example: Overflow, Underflow, Division by Zero, Illegal memory address, Mis-aligned memory access

Synchronous

- Occur due to program execution
- Have to execute an **exception handler**

Asynchronous

- Occur **independently** of program execution
- Have to execute an **interrupt handler**

Threads

- A process may have multiple independent control flows called threads
- Each thread has its own stack and registers (PC, SP, registers), but share the same address space
- Shared memory model and Shared memory architecture
- Faster thread generation- no copy of address space
- Different process can be assigned to run on different cores of a multicore processor

User threads

- Managed by library
- Context switch is fast, OS not involved
- **Disadvantage**

1. OS cannot map different threads of the same process to different resources \Rightarrow No parallelism
2. OS cannot switch to another thread if one thread blocks

Kernel threads

- OS is aware of the threads and can manage accordingly
- Efficient in a multicore system
- Potential synchronisation issues

Many to one mapping

- All user-level threads mapped to one process.
- Efficiency depends on threading library

One to one mapping

- Each user-level thread is mapped to one kernel thread
- OS schedules

Many to many mapping

- Many user-level threads mapped to many kernel threads
- Library threads has overheads, and kernel threads has overheads
- At different points in time, different user threads are mapped to different kernel threads
- Number of threads must be suitable to the degree of parallelism and the resources available

Locks

- Spinlock: busy wait
- Blocking: mutex
- Using more locks increases the number of context switches
- DO NOT wait in the critical section

Semaphores

- Essentially shared global variables
- Can be potentially accessed anywhere in program
- No connection between semaphore and the data being protected

Barrier

- All threads must reach the barrier before any thread can proceed

Deadlock

- Deadlock exists among a set of processes if every process is waiting for an event that can be caused only by another process in the set
- **iff these condns are met**

1. Mutual exclusion-at least one resource is not shareable
 2. Hold and wait - at least one process holding a resource and waiting for another
 3. No preemption - critical section cannot be aborted externally
 4. Circular wait
- Dealing with deadlock**
- Ignore it, prevent it, avoid it by controlling resource allocation, detection and recovery by breaking cycles

Starvation

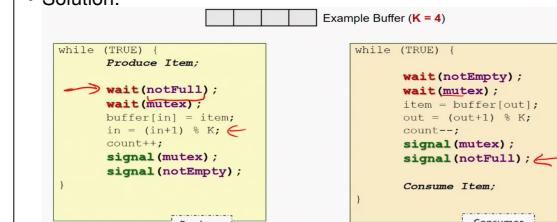
- Side effect of the scheduling algorithm. Lower priority processes might starve

Livelock

- Active acquire release but no useful work done

Producer-Consumer Problem

- Specifications:
 - Producers put in a shared bounded buffer if not full, consumers read from it if not empty
- Solution:



- Concurrent read, exclusive write. Categorical starvation of writer is possible

- | Writers | Readers |
|--|---|
| <pre>roomEmpty.wait ()
#critical section for writers
roomEmpty.signal ()</pre> | <pre>mutex.wait ()
readers += 1
if readers == 1:
 roomEmpty.wait () # first in locks
 mutex.signal ()
critical section for readers</pre> |
| | <pre>readers -= 1
if readers == 0:
 roomEmpty.signal () # last out unlocks</pre> |
| | <pre>Light switch: Abstracts out the shared lock for the reader
counter = 0
mutex = Semaphore (1)
lock (semaphore):
 mutex.wait ()
 counter += 1
 if counter == 1:
 semaphore.wait ()
 mutex.signal ()
unlock (semaphore):
 mutex.wait ()
 counter -= 1
 if counter == 0:
 semaphore.signal ()
 mutex.signal ()</pre> |
| | <pre>Starvation free solution (block out readers):</pre> |

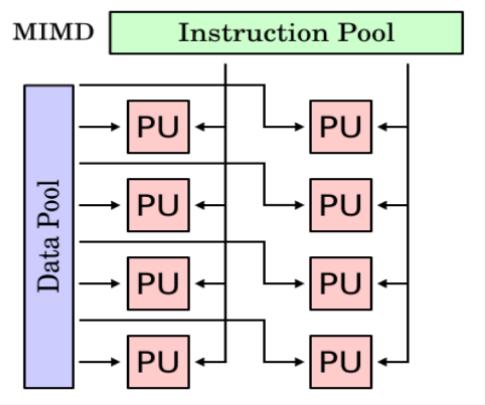
- Same instruction broadcasted to all ALUs
- AVX: intrinsic functions operate on vectors of 4 64 bit values

Multiple Instruction Single Data

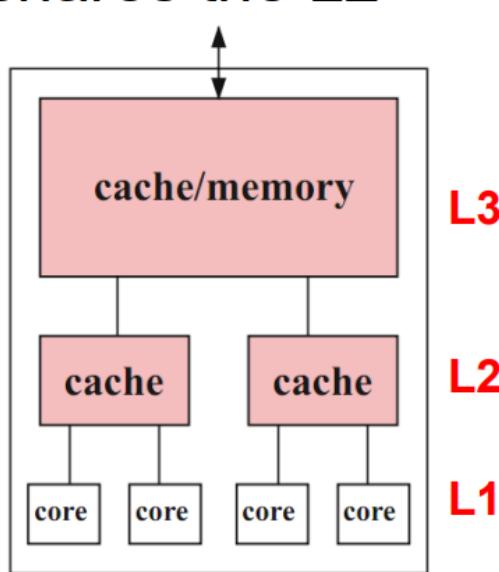
- Multiple instructions operating with a single data

Multiple Instructions Multiple Data

- Each PU fetches its own instructions
- Each PU operates its own data
-

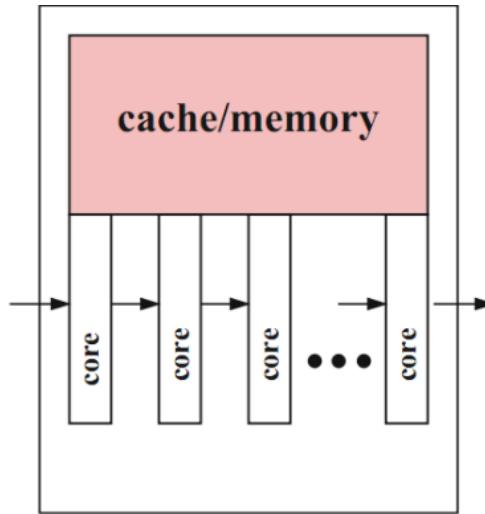


Hierarchical designs



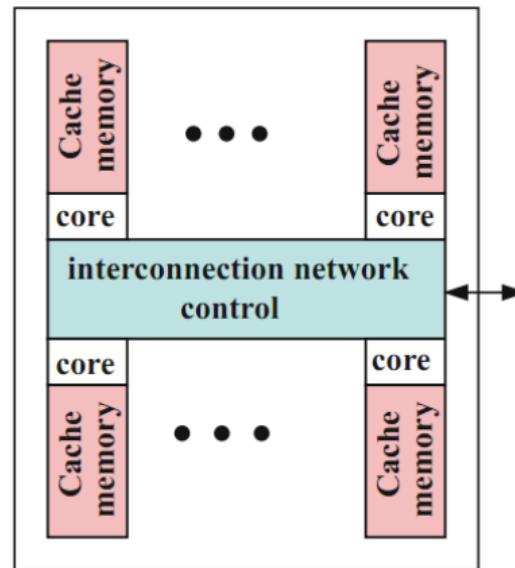
- Each core can have a separate L1 cache and shares the L2 cache
- All cores share common external memory

Pipelined design



- Multiple packets being processed in a pipelined fashion
- Cores connected linearly, shares the same cache, memory
- Useful if the same computation has to be applied to a long sequence of data elements

Network-based design



- Cores and their local memory and memories are connected via an interconnection network

Why cache

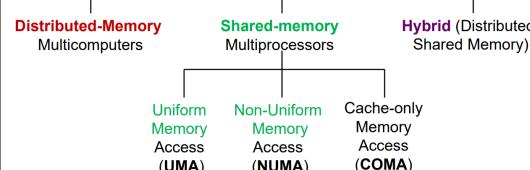
- Cache provides high bandwidth data transfer to CPU and reduce latency in data access
- Memory latency: Amount of time for a memory request from a processor to be serviced
- Bandwidth: Rate at which the memory system can provide data to a processor

- A stall happens when the next instruction depends on previous instructions
- Bandwidth and latency affects stalls, since instructions (sw, lw) needs to wait for the memory system to become available

Performant parallel programs

- Try not to overload the memory system with too many requests
- Share data across threads (inter-thread cooperation)
- Reuse data fetched previously (temporal locality)
- Favor additional arithmetic over load / store

Parallel Computers



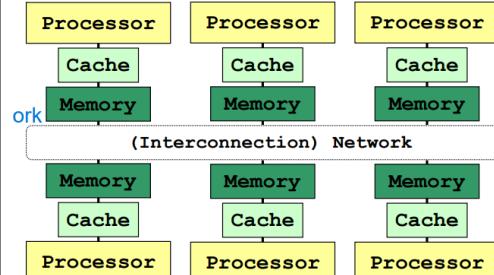
Cache coherence

- Multiple copies of data exist on different caches
- Local updates should not be seen by other processes
- Maintained by additional instructions
- Instructions that mess up cache coherence hence presents severe overheads

Memory consistency

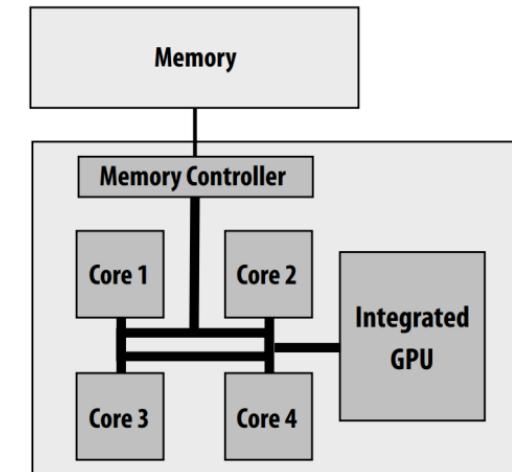
- Memory consistency depends on the PL and architecture
- A seq consistent architecture makes a PL with seq const memory model run faster since fewer instructions are needed to ensure memory consistency

Distributed Memory



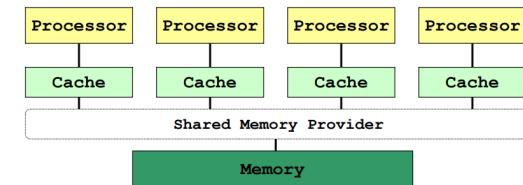
- Each node is an independent unit with processor and memory
- Memory in each node is private
- Nodes communicate through a network

Shared memory



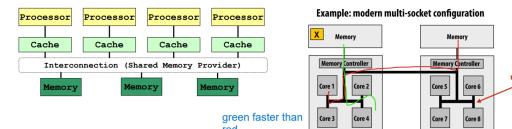
Intel Core i7 (quad core) (interconnect is a ring)

- Parallel programmes share memory through controller / provider
- Cache coherence and memory consistency is ensured



Uniform Memory Access

- Latency of accessing main memory is the same for processors
- Suitable for small number of processors. Contention over memory can be high for large number of processes

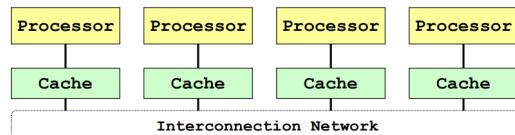


Non-uniform Memory Access

- Physically distributed memory of all processing elements are combined to form a global shared memory
- Local memory access has lower latency
- Reduce contention since each processor tends to access local memory
- Adding more processes does not increase contention as much as UMA
- Data consistency is easier too

Cache Coherent NUMA (CCNUMA)

- Each node has cache to reduce contention



Cache only Memory Architecture (COMA)

- Each memory block works as cache memory. This means that no fixed space stores data permanently and cache block with data can be moved around dynamically.
- Data migrates dynamically to keep data as close as possible to the processors
- Cache coherence is harder since data may not just be copied, they can also be shifted around.

L7: Cache coherence and memory consistency

L11: Interconnection networks

Parallel Computation Models

L4: Shared-memory programming models

Parallelism

- Average number of units of work that can be performed in parallel per unit time.
- E.g. MIPS, MFLOPS
- Limitation: Program dependencies - data, control
- Runtime delays - memory contention, communication overheads, thread overhead, synchronisation
- We cannot reorder them however we like
- Work = Task + dependencies (limitations)

Data parallelism

- If iterations are **independent**, they can be executed in arbitrary order on multiple cores
- Partition data among processing units, each doing similar work
- Commonly expressed as a loop, if the iterations are independent and can be executed in arbitrary order
- E.g. SIMD computers
- OpenMP - matrix multiplication**

```
// parallelize result = a * b
// each thread works on one iteration of
// the outer-most loop
// vars (a, b ,result) are shared
#pragma omp parallel for num_thread(8)
    shared(a, b, result) private(i, j ,k)
    ...

```

Same as

```
for (i=0; i < size; i++)
    for (j=0; j < size; j++)
        for (k=0; k < size , k++)
            result[element][i][j] +=
                a.element[i][k] *
                b.element[k][j]
```

Single Program Multiple Data (SPMD)

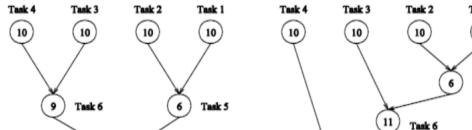
- Same programme may behave differently based on the data
- E.g. Scalar product of $x.y$ on p processing units

```
local_size = size/p;
local_lower = me * local_size;
local_upper = (me+1) * local_size - 1;
local_sum = 0.0;

for (i=local_lower; i<=local_upper; i++)
    local_sum += x[i] * y[i];

Reduce(&local_sum, &global_sum, 0, SUM);
```

Same program executed by p processing units.
"me" is the processing units index (0 to p-1)

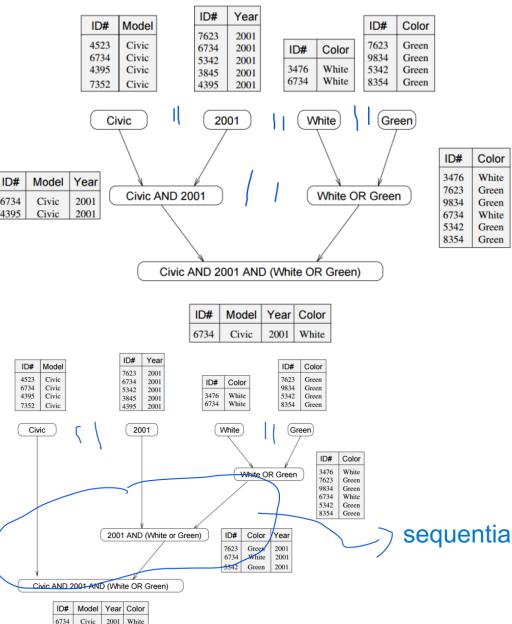


Critical Path = (Task 4 → 6 → 7)
Critical Path Length = 27
Degree of concurrency = 63 / 27 = 2.33

Critical Path = (Task 1 → 5 → 6 → 7)
Critical Path Length = 34
Degree of concurrency = 64 / 34 = 1.88

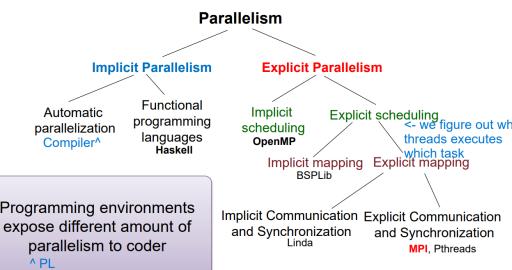
Task parallelism

- Partition the tasks among the processing units
- independent program tasks/ parts can be executed in parallel
- Granularity: statement, loop, function
- More complexed than data parallelism → needs to schedule, map, take care of dependencies ...
- Decomposition**
 - The room for parallelism in a task depends on how the task is decomposed



Task dependence graph

- DAG: node=tasks, value=expected execution time, edge=control dependency
- Bad for one process to take disproportionately more data → idle time
- Critical path length: maximum slowest completion time
- Degree of concurrency=total work/critical path length



Programming environments expose different amount of parallelism to coder ^ PL

Coordination: Shared memory

- Protect access to shared address space, mutex.
- Needs hardware support to implement efficiency. NUMA makes it easier but it is still costly to scale due to contention (any processor can load/ store to any address)
- Can be done without a shared memory system (NUMA, UMA)
- Any type of coordination can be used in any hardware via software

Coordination: Data-parallel

- SIMD, vector processors
- Traditional: Map a function onto a large collection of data
- Side effect free execution
- Modern: Data-parallel languages do not enforce this structure
- SPMD model used in CUDA, OpenCL, ISPC instead

Coordination: Message passing

- Tasks operate within their own private address space and communicate by explicitly sending / receiving messages
- E.g. MPI, GO
- Hardware does not implement system wide loads and stores, can connect commodity systems together to form large parallel machines
- Many many computers, not a very big one
- Compatible with distributed memory systems

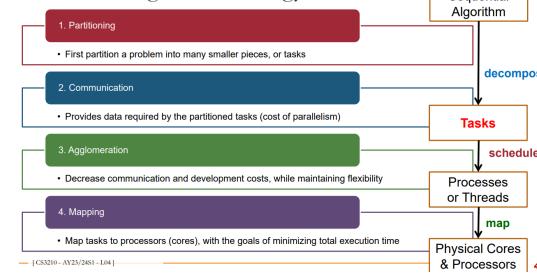
Coordination and hardware

- Shared memory: UMA, NUMA. Copies of messages and sent / received from library buffers
- Message passing: distributed systems, clusters, supercomputers
- Any abstraction can be implemented with any hardware but it will be more costly
- Shared address space on incompatible hardware
 - Write: Send message to all cores to invalidate value
 - Read: page fault handler issues appropriate network requests

Summary of Coordination Models

- Shared address space: very little structure
 - All threads can read and write to all shared variables
 - Drawback: not all reads and writes have the same cost (and that cost is not apparent in program text)
- Data-parallel: very rigid computation structure
 - Programs perform the same function on different data elements in a collection
- Message passing: highly structured communication
 - All communication occurs in the form of messages

Foster's Design Methodology



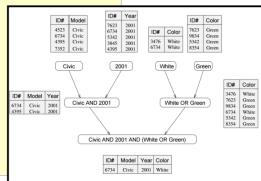
Foster's Design methodology

- Partitioning**
 - Divide computation and data into independent pieces to discover maximum parallelism
 - Two approaches:**
 - Domain decomposition: divide data into smaller, equal pieces. Associate computation with data.
 - E.g. 24 tasks with 3 grids each → 6 tasks with 12 grids each
 - Functional decomposition: Divide computation into piece. Associate data with computation.
 - E.g. Climate model → Atmospheric model, hydrology model ...
- Rule of thumb:**
 - 10x more primitive tasks than cores in target computer
 - Minimize redundant computations and redundant data storage
 - Primitive data should be of roughly the same size
 - Number of tasks an increasing function of problem size
- Communication (coordination)**
 - Dependencies between tasks necessitates communication
 - Overlap computation and communication such that when some tasks are communicating, others are computing (improve utilisation)
- Local Communication**
 - Tasks need data from a small number of other tasks (neighbors)
 - Use channel
- Global Communication**
 - Significant number of tasks contribute to perform a computation
 - Do not create channels early on in the execution
- Parallel Programming Patterns**
 - Pattens are not mutually exclusive, use the best match
- Fork Join**
 - Children run in parallel but are independent
 - Children execute the same or different program

- Children join the parent at different points
- Implementation:** Processes, threads etc

```
P1 = Fork {
    P3 = Fork { return Model = "civic" }
    P4 = Fork { return Year = "2001" }
    Join P3, P4
    Return P3 AND P4
}

P2 = Fork {
    P5 = Fork { return Color = "green" }
    P6 = Fork { return Color = "white" }
    Join P5, P6
    Return P5 OR P6
}
Join P1, P2
Return P1 AND P2
```



Parbegin - Parenend

- most relaxed, code is structured into sequential segments and parallel segments
- Programmer specifies a sequence of statements to be executed in parallel
- A set of threads is created and the statement of the construct are assigned to these threads
- All the forks are done at the same time and all the joins are done at the same time
- Statements after parbegin and parenend are only executed

after all threads joins (barrier)

- Implementation:** OpenMP or compiler directives
- E.g Matrix multiplication using openMD

SIMD (not the Architecture)

- Single instructions are executed synchronously by different threads on different data
- Similar to parbegin-parend but all threads execute the same instruction at the same time (synchronous)
- Parallel but synchronous
- Implementation:** AVX / SSE instruction on intel processor

```
mulps xmml, xmm0
127 95 63 31 0
4.0 3.0 2.0 1.0
* * * *
5.0 5.0 5.0 5.0
= = = =
20.0 15.0 10.0 5.0
```

xmm registers are 128 bits long

SSE instruction treats the xmm registers as 4 individual 32-bit floating point value

SPMD

- Same program executed on different cores but operate on different data
- Different threads might execute on different instructions of the same program due to control flow (ifs) and speed of cores

- Similar to parbegin-parend but there is no implicit synchronization (lack of barrier)
- E.g. programs on GPGPU

Master-Worker

- Single program controls the execution of the program
- Master executes main function, assigns work to worker threads
- Initialisation, output and Coordination is done by master
- Worker waits for instruction

L6: Data parallel models (GPGPU)

L9,10: Distributed-programming models

Performance and Scalability of Parallel Programs

L5: Performance of parallel systems

Two Views

- Response Time (user): duration of a program is reduced (start - end time)
- Throughput (computer manager): more work to be done in the same time (jobs per second)

Performance Factors

- Programming Model
- Computational Mode: How well the given program runs in the given architecture
- Architectural Model: interconnection network, memory organization, execution mode, sync or async processing

Response time in sequential programs

- Wall-clock time
- Comprise of
 - User CPU time: time CPU spends executing program
 - Know that read and write cycles take different time
 - System CPU time: time CPU spends on system instructions. Depends on OS.
 - Waiting time: IO waiting time and execution of other programs due to time sharing. Depends on the load of the system.

L8: performance instrumentation

New Trends

L12: Energy efficient computing