

Introduction to Scalable Systems

Data Structures & Algorithms

→ Saving storage for sparse matrices

→ Each non-zero element is added as a tuple in a list.

→ This tuple contains the row, coln., & value.

→ Storing row & coln. overheads is better than storing all 0s.

→ Compressed Sparse Row (CSR) representation (Stores as three arrays).

→ $V[\text{nnz}]$: non-zero values in row-major order

→ $C[\text{nnz}]$: coln. offset within row group for a non-zero value

$\text{no. of non-zero elements} \rightarrow R[M+1]$: stores cumulative count of non-zero elements till $(i-1)^{\text{th}}$ row
dimension of the matrix

→ Relatively cache friendly

Hashing Analysis

→ Using chain hashing: Expected no. of keys in successful search

$$= \frac{1}{n} \sum_{j=1}^n \left[1 + E\left(\sum_{i=i+1}^n X_{ij}\right) \right]$$

$$= 1 + \frac{n-1}{2b}$$

$$= O(1+\alpha), \text{ where } \alpha = \frac{n}{b}$$

B-Trees

$$\rightarrow \text{Depth of the tree} = \left\lceil \log_{\frac{m}{2}} (\text{size}) \right\rceil + 1$$

Graphs

→ Diameter = Distance of the longest shortest path!
 $= \max d(u, v) \quad \forall u, v \in V, u \neq v$

→ Clique : A subgraph which is a connected graph
 → Maximal clique : Clique with most no. of nodes possible.

→ Graph traversals
 ↴ BFS

→ Start at a source vertex. Use a queue

→ $O(|V|)$ for adjacency matrix

→ $O(d)$ for " list

better for sparse graph

better for dense graphs due to cache locality

→ Total time

$w = \text{no. of vertices in connected component}$

$\rightarrow O(w|V|)$ for adjacency list

$\rightarrow O(w + f)$ " " matrix

no. of edges in connected component.

→ DFS (Use stack instead of queue)

→ Requires $O(1)$ space (for recursion stack)

→ Same TC as BFS

Algorithms

→ Dijkstra's algorithm

Worse than array implementation

When $|E| \gg |V|$, i.e. for

an almost-connected graph

Method of storage	TC
Array of weight	$O(V ^2 + E)$
Using min. heap to store edge weights	$O((V + E) \log V)$
Using Fibonacci heap	$O(E + V \log V)$

Parallel Programming

→ Flynn's classification : In terms of parallelism & data stream

→ SISD: Single instruction single data (parallel programming)

→ SIMD: " " multiple " (vector processor & processor arrays)

→ MIMD: Multiple " " "

→ Used in supercomputers etc.

→ Classification based on memory

→ Shared memory : UMA & NUMA

→ MVMA (Multi-Uniform memory architecture)

→ Message passing

→ UMA (Uniform memory architecture)

→ Shared memory is distributed among processor nodes

Shared Memory & Cache

→ If a mem. location has been updated by a processor, how will other processors know? They may use the old value (in their cache)!

→ Solutions:

→ Write update: Propagate cache line to all processors on write write.

→ Writes shall be coherency

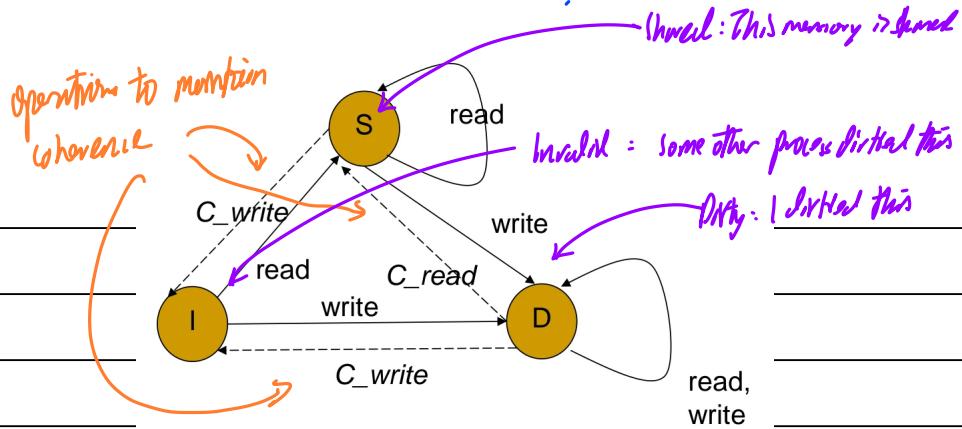
→ Cache line is sent to all processors, even those who do not need it, thereby wasting resources.

Followed in all systems →

→ Write invalidate: Inform all other processors that they have "stale" data, and to read from memory next time.

→ Read via writer, dirty bit is valid.

→ Three states are possible:



Implementation

- Snoopy: Used with bus-based architecture
- Directory based:
 - A small part of shared memory is used to cache states, so that any changes can be relayed to relevant processors only.
 - Contains cache states and the list of relevant processors for a cache line.
 - Increases memory access time.
 - Implementation:
 - Use presence bits for owner processors:
 - Full-bit vector scheme: $O(M \times P)$ storage for P processors, M cache lines. (Not required)
 - Sparse (try-and) directory: limited cache lines & presence bits
- False sharing
 - A cache block may be shared by two processors, but not the actual memory location.
 - Solution: restrictive code s.t. it is non-intrusive
 - No pushing to ensure a cache line only contains own data.

Interconnect

- PCI or PCI-e is used to connect a processor to a network line.
- Communication network: consists of switching elements connected w/ each other
- Topologies:
 - Bus / ring
 - Crossbar switch: 2-D grid of switching elements
 - P^2 switching elements

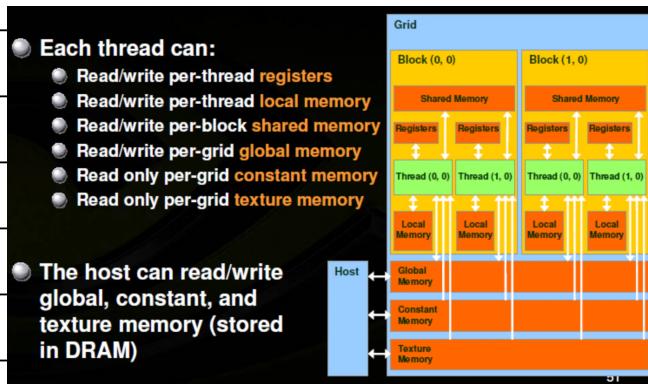
- Multi-stage network (omega)
 - consists of $(\log_2 P)$ switches, each with $P/2$ switching elements
 - can be blocking
- Mesh, torus, hypercube, and fat-tree
- Evaluation metrics for interconnection topologies:
 - min. no. of links to be removed to disconnect the topology → Diameter: max dist. b/w any two nodes
 - Connectivity: multiplicity of paths b/w 2 nodes
 - Bisection: min. no. of links to be removed to split the network to two equal halves.
 - Channel width: no. of bits that can be simultaneously communicated over a link.
 - " rate: perfermance of a single physical wire
 - " bandwidth: channel width \times channel rate
 - Bisection " : min. volume of communication b/w two halves of the network
 - bisection width \times channel bandwidth

Type of network	Diameter	Type of Network	Connectivity
Full-connected	1	Linear array	1
Star	2	Ring	2
Ring	$P/2$	2-D mesh	2
Hypercube	$\log_2 P$	2-D mesh (w/ wraparound)	4
		D-dimension hypercube	D

Type of network	Bisection width
Ring	2
P-node 2-D mesh	\sqrt{P}
Tree	1
Star	1
Completely connected	$P^2/4$
Hypercubes	$P/2$

Memory Hierarchy

- Global/device memory
 - Can be accessed by all threads in the streaming multiprocessor(SM)
- Shared memory
 - In each SMX, and shared by all threads in a SMX.
- Can be accessed by CPU (host)
- 200-400 clock cycles per access
- Can be configured as a 1KB shared mem., and the rest as L1 cache
- 20-40 clock cycles for access



- Differences w/ CPU threads:
 - Fast context switching: A thread's context is stored in shared memory (zero overhead context switching) & registers until it completes.
 - Explicit cache management: User's program will have to bury the frequently used data manually.

Parallelisation Principles

- Overheads
- Communication delay → Illing → Synchronisation
- To prevent this, ensure program is load balanced

Evaluation metrics

$$\rightarrow \text{Speedup}, S(p, n) = \frac{T(1, n)}{T(p, n)} \quad // \text{Time for sequential execution}$$

$T(1, n)$ // Time for parallel execution with p processors.

$$\rightarrow \text{Ideally}, S(p, n) = p \quad // \text{Equal distribution of loop}$$

- Usually, $S(p, n) < p$ // due to overheads
- Rarely, $S(p, n) = p$ // where we get added benefits due to parallelisation, i.e. cache

$$\rightarrow \text{Efficiency, } E(p, n) = \frac{S(p, n)}{p}$$

- Speedup limits: Amdahl's law // law of diminishing return

- Performance is always bottlenecked by sequential section
- Places a limit on speedup:

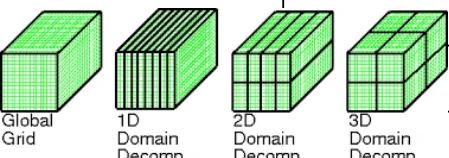
$$S(p, n) \leq \frac{1}{f_s + \frac{f_p}{p}} \quad \begin{matrix} \\ \text{|| } f_s = \text{time for sequential} \\ \text{section} \end{matrix}$$

$f_p = \text{time for parallel section}$

- Gustavson's Law: Increase the problem proportionally to keep the overall time constant.

→ Parallelizing a program

- Step 1: Decomposition & Step 2: Assignment



- Identify segments of code that can be parallelised.

- How to group tasks? // use structured approach

- Roman decomposition

- Divide the data among nodes, each process owns & computes a portion of the data

- This is specified by a process grid.

- Data distribution: Schemes used to divide data

- Blocks: for regular (uniform) compute

- Block-cyclic: for imbalanced compute. Divides the data into blocks and distributes them among the processes in a cyclic fashion.

- Step 3: Orchestration: Reducing synchronization overhead

- Structures & synchronizers communication

can be done by packing data structure \hookrightarrow

→ goals :

→ Maximise data locality

→ Minimising volume of data exchange

→ Minimise frequency of data exchange

→ Reduce communication by replicating data

E.g. not communicating intermediate results

→ Minimise contention / hotspots

→ Do not use same communication pattern as the other processes.

→ Overlap iterations with computations

→ Sacrifice compute to save on communication.

→ Step 4: Mapping

→ Which process runs on which processor?

→ Dynamic Mapping / self-scheduling / work-stealing

→ A process holds a set of tasks, which it distributes to nodes as and when they become free.

OpenMP

→ Shared memory using compiler directives

→ Uses fine-grain model; provides fine-grained parallelism

→ Supports nested parallelism

→ Constructs

→ `parallel` : #pragma omp parallel < clauses >

→ `private (list)` : Declares the vars. in list to be private to each thread

→ `shared (list)` : Declares the vars. in list as shared by all threads

→ `firstprivate (list)` : `private (list) + initialisation`

→ `default (shared|none)` : Sets default data copying rule if variables within the parallel region.

→ `#pragma omp if(exp)`: Executes the section in parallel if expression is true. Otherwise, the segment is executed in sequential.

→ `#pragma omp num_threads(int)`: Defines a number of threads to use for the parallel region.

→ `#pragma omp reduction(Op list)`: Reduces the values from all threads into a single variable.

→ For: # pragma omp for < clause>

→ schedule clause:

→ `schedule(static, chunk-size)`: original data is split in chunks and divided among the threads in

→ `schedule(dynamic, chunk-size)`: Original data are distributed among the threads in a dynamic manner.

→ `schedule(runtime)`: chunk size, no. of threads, etc. are decided at runtime.

→ Synchronisation Directives:

→ `#pragma omp master`: this section can be executed only by master

→ `#pragma critical`: defines the critical section

→ `#pragma atomic`: this section must be executed w/o context switching

→ `#pragma flush(var-list)`: flush all variables on var-list set. the update's made to those variables are visible to all threads

→ `#pragma ordered`: Threads should execute in order of thread no.

Message Passing Interface - MPI

→ Point-to-point communication:

* Rank is essentially the unique id of a process w.r.t. the particular communication context. It nodes can have different ranks in different groups.

MPI-SEND (buf, count, datatype, dest, tag, comm)

- Message
- Destination rank
- Communication context: used to indicate the group where the message must be sent

MPI-RECV (buf, count, datatype, source, tag, comm, status)

- Message
- Source rank
- Sender & receiver rank
- Communication context: receiver tag should match tag given by MPI initiator

* MPI_COMM_WORLD is the wild card communication context. To check the MPI_Comm_rank is used by a process to obtain its own rank.

Receivers source & tag field can be wildcarded by using MPI_ANY_SOURCE & MPI_ANY_TAG respectively.

→ Some utility functions:

→ MPI_Init: Initialize MPI environment

→ MPI_Finalize: (close)

→ MPI_Comm_size (comm, &size): Used to get total no. of processes in the group pointed to by the communicator

→ MPI_Wtime: Returns current wall clock time.

→ Non-blocking communication

stands for immediate because these functions return immediately

→ MPI_ISEND (buf, count, datatype, dest, tag, comm, request)

→ MPI_IRECV (buf, count, datatype, source, tag, comm, request)

MPI-WAIT (request, status)

A sort of token for checking request status

blooms until communication is complete

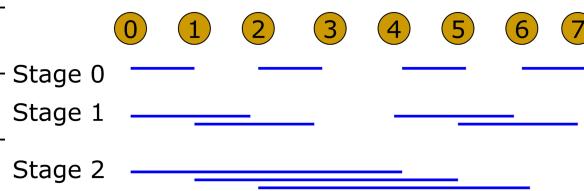
MPI-TEST (request, flag, status) → checks if communication has completed

MPI-REQUEST-FREE (request) → used to de-allocate (force stop) a process.

→ Collective communication using barrier will result in blocking, because a thread is waiting for all other threads. To solve this, we can use the butterfly algorithm!

→ In the n^{th} round, each process i synchronizes with process $i \oplus 2^k$ pairwise.

→ Requires $2 \log P$ piecewise synchronization



CUDA Programming

→ A kernel is executed by a grid of thread blocks. A thread block is a batch of threads that can co-operate with each other by:

→ Sharing data using shared memory, or

→ Synchronising their execution

→ Threads from different blocks cannot co-operate

→ CUDA Memory spaces:

Memory	locked?	Access Modes	Scope
Local	X	RW	One thread
Shared	N/A	RW	All threads in a block
Global	X	RW	All threads & host
Constant	✓	R	
Texture	✓	R	

- Differences b/w LPV & CUDA threads
 - CUDA threads are lightweight
 - zero-overhead context switching
 - very minimal creation overhead
 - (CUDA vs. 1000s of threads)

→ CUDA is a minimal extension of C

→ CUDA Built-in variables:

- $\text{threadIdx.}\{x, y, z\}$: thread ID within a block
- $\text{blockIdx.}\{x, y, z\}$: block " " " grid
- $\text{blockDim.}\{x, y, z\}$: no. of threads in a block
- $\text{gridDim.}\{x, y, z\}$: " " blocks " " grid

→ CUDA function

- `cudaMalloc & cudaFree`: Used to allocate & free memory on GPU
- `cudaMemcpy`: to copy data from LPV to GPU
- `syncthreads()`: to sync threads in a block (Lorrier)

PRAM

→ Ideal model

→ Unlimited processors → Shared global memory

→ Uniform access time → SIMD execution

→ Handling conflicts

→ ER EW: exclusive read, exclusive write

→ REW: common " " " (standard)

→ CRCW: " " common "

→ Common: All processors are writing the same value

→ Arbitrary: One random process is allowed

→ Priority: Processor with least PID is allowed

→ Reduction: Values are reduced (by sum/avg)

Parallel Quick Sort

- One processor identifies the pivot & relays it to all processors.
- Each processor does local partitioning
- Processors are grouped based on the partitions to ensure load balancing
- All elements lower than the pivot will be shifted to group 1, and all other elements go to group 2. We use prefix sums to find the correct location of the data.
- Eventually, P_0 will get the smallest numbers, P_1 will get next set of smallest numbers, etc. Each processor does a regular quick sort and then the data is merged.
- Complexity: $\underbrace{O\left(\frac{N}{P} \log\left(\frac{N}{P}\right)\right)}_{\text{communication}} + \underbrace{O\left(\log^2 P + \frac{N}{P} \log P\right)}_{\text{compute}}$