

Unit 2

lec 10-20

PRAm is gud cuz its easy to design algorithms, analyse.

Work Time Scheduling Principle

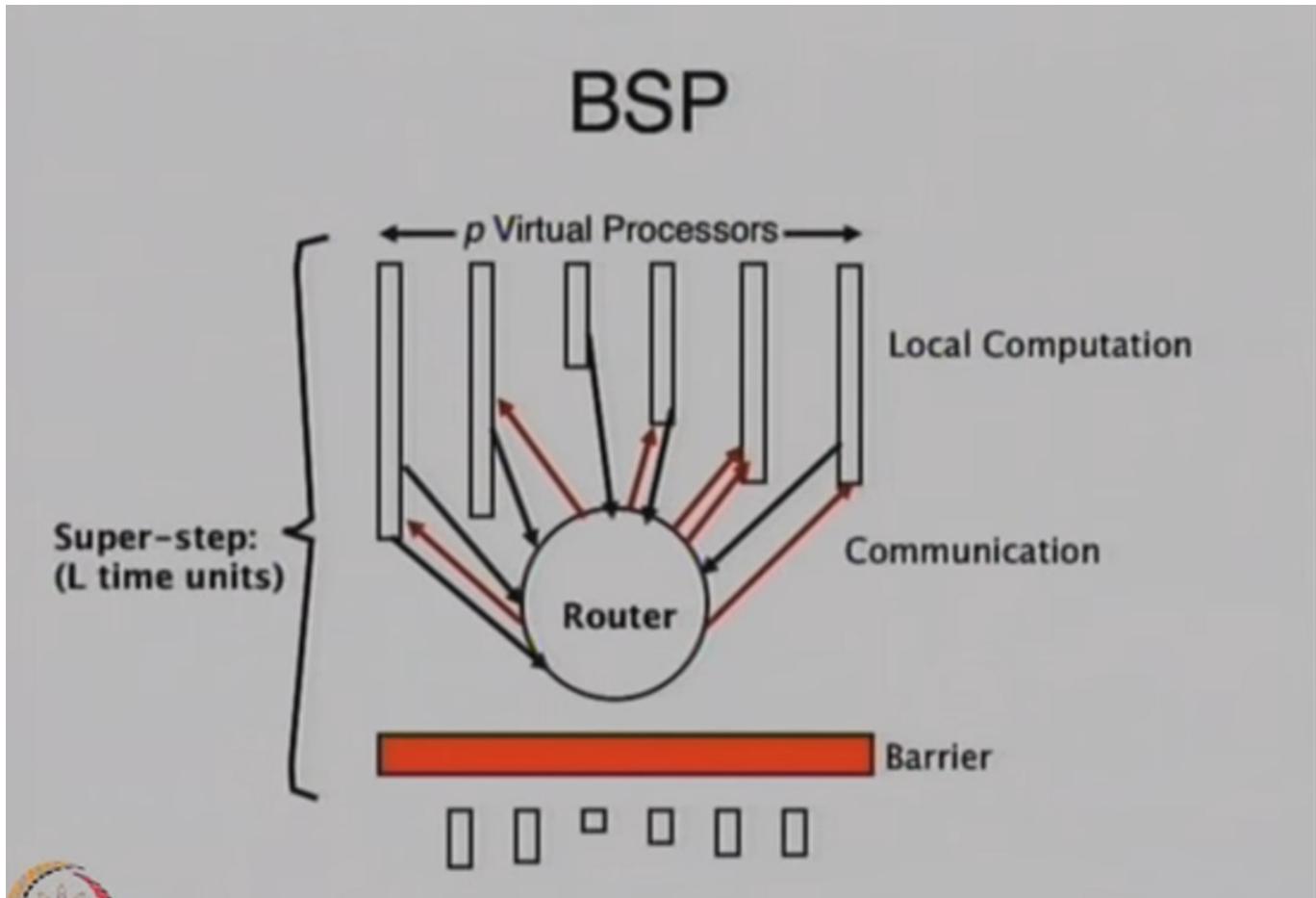
This essentially says do not care about the actual number of processors present. Make a algorithm that evaluates in $t(n)$ steps using p processors and let a auto scheduler manage the algorithm in the bare metal.

Cost optimal algorithm is automatically work optimal (this is compared to sequential).

Bulk Synchronous Parallel Model

This says that there are virtual pairs of processor and memory. There is not local mapping but instead network memory. It operates by executing "super step":

- local computation
- communication
- barrier synchronization



BSP has system wide steps, like pram. It is simple to analyze. It connects well to physical architecture.

Calculating cost for this is a little more complex.

- barrier cost of l (barrier sync)
- superstep cost: local computation + communication + barrier
 $\Rightarrow w + hg + l$ g is thruput, h is message something
- $Total/cost = \sum w_s + g \sum h_s + Sl$

LogP Model

btw this isn't math log. This is l,g,p from above. wat is o? no one knows.

Let's say we are doing a regular search. N elements and p processors. This will be technically faster than regular binary search **but** WORK DONE is way higher.

Memory Consistency

Cache Coherency

This is L1 cache of a thread making sure to write to L2, or whatever is the shared memory, to maintain that data's "freshness" or coherency. This makes sure that all threads work with proper updated data.

Snoopy cache coherence protocol is literally wiretapping.

Strict Consistency

- `read(x)` must return the latest `write(x)`. It should also appear as if operations are instantaneous.
- linearizability is a weaker version. It has a point for sync called linearization point. It is instantaneous here.
- Impractical: wth is this? It said it is done for uniuser, so its just nothing happening? o_o

Sequentially Consistent Model

This is where the memory is consistent to how the memory would have been if the program was sequential. This is hard to implement.

Processor Consistency

All processes must see memory writes from one process in the order they were issued. This is FIFO consistency. Processor consistency is when we say all processes see memory writes to a variable (data) must be seen in the order of issue.

Weak Consistency

Consistency is only enforced on request. OpenMP uses this.

Summary of Consistency Models

Model	Description
Strict	Global time based atomic ordering of <i>all</i> shared accesses
Sequential	All threads see all shared accesses in the same order consistent with program order -- no centralized ordering
Causal	All threads see causally-related shared accesses in the same order
Processor	All threads see writes from each other in the order they were made. Writes from different processes may not always be seen in that order
Weak	Special synchronization based reordering -- shared data consistent only after synchronization

Performance issues

- True sharing:
 - make copies for each processor
- False sharing
 - give continuous blocks . These blocks have the variables. In this we send lines of memory. So we get some useless stuff.

There are other things like make code short, use as many threads as required etc.

Message Passing Interface (MPI):

This is just what the heading is. No magic nonsense. Also the low level implementation is left out.

There is synchronous or asynchronous implementations.

In async we are using buffers and stuff to CHECK sync. Sync enforces the CHECK between sender and receiver.

- Standard mode: This will choose from the following modes dynamically.
- Buffered mode: I dunno. User specifies the buffer.
- Synchronous mode: Will complete once receive operation is accepted. We can send stuff without waiting for a corresponding receive.
- Ready mode: Sending waits for a receive. This is like standard but with an optimization

Messages are:

- in-order
- have progress (send or receive will complete)
- not fairness guaranteed (receive can be open to all or open to one)
- resource limited, sync are usually more resource efficient

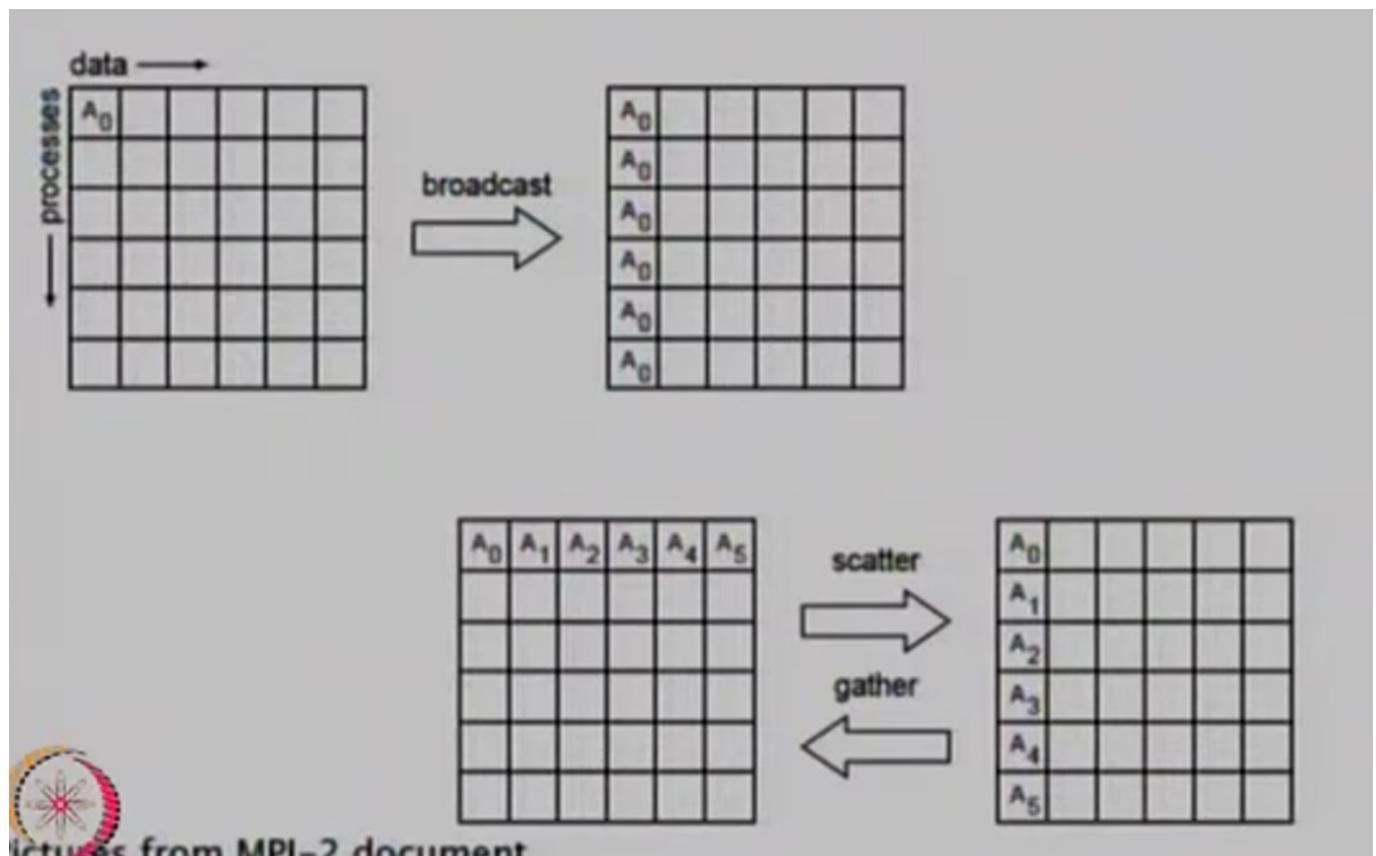
They have data in a special format, `MPI_INT`, `MPI_CHAR` ... etc as well as user defined one like `MPI_<whatever>` and we can make `MPI_VECTOR` which is like blocks of data.

For user defined types we have to "commit" the datatype, after creating the type, before we use it in communications. This is an OpenMPI specific thing. Other don't use "commits".

There are also a few different types of "collective" or group communications.

1. `MPI_Barrier`: This is like a message that causes a sync. It pauses threads until all (or the required threads) of them have sent this message. Imagine like waiting for your entire team to finish the lap before all of you can go to the next lap. This only returns after everyone calls Barrier.

2. MPI_BCast: One message to all. Broadcast. This returns after YOU receive the message. No guarantee of the others receive. Broadcast has to have previously chosen sender and receiver, and both make a call to send and receive. It is not interrupt based. A broadcast happens to a communicator (basically a yet channel). A blocking broadcast REQUIRES everyone to read before other broadcasts.
3. MPI_Scatter: One unique piece of the message to each one.
4. MPI_Gather: Get all parts from a group.
5. MPI_Alltoall: All Send, All Gather.
6. MPI_Reduce: This is reduce from the threads sync i.e. take all results and compile according to the requirement. |||ly MPI_Reduce_Scatter.
7. MPI_Scan: Not what it seems. Get all values before you, process it and do something ig. Kinda weird tbh. Google for clarity or more pain.



A ₀					
B ₀					
C ₀					
D ₀					
E ₀					
F ₀					

A ₀	A ₁	A ₂	A ₃	A ₄	A ₅
B ₀	B ₁	B ₂	B ₃	B ₄	B ₅
C ₀	C ₁	C ₂	C ₃	C ₄	C ₅
D ₀	D ₁	D ₂	D ₃	D ₄	D ₅
E ₀	E ₁	E ₂	E ₃	E ₄	E ₅
F ₀	F ₁	F ₂	F ₃	F ₄	F ₅

In-place MPI_Scan

process i receives data reduced on process 0 to i

	recvbuf						recvbuf				
P0	3	4	2	8	1		3	4	2	8	1
P1	5	2	5	1	7		8	6	7	9	8
P2	2	4	4	10	4		10	10	11	19	12
P3	1	6	9	3	1		11	16	12	22	13

```
MPI_Scan(MPI_IN_PLACE, recvbuf, 5, MPI_INT, MPI_SUM,
          MPI_COMM_WORLD)
```

We can "spawn" threads with `MPI_Comm_Spawn`. Window is like a chuck of memory I suppose. Use `MPI_Put` and `MPI_Get` to work with it.

There are so many other ones man.

Remote Memory Synchronization

- MPI_Win_fence
- MPI_Win_lock
- MPI_Win_unlock
- MPI_Win_start
- MPI_Win_complete
- MPI_Win_post
- MPI_Win_Wait
- MPI_Win_Test

Algorithm Techniques

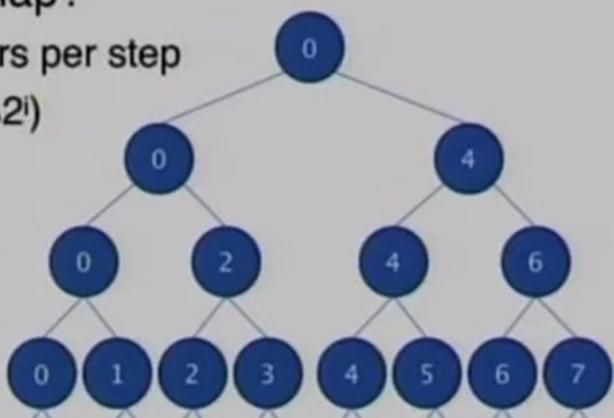
Balanced tree approach:

We build a binary tree on the input, thereby dividing into hierarchical groups. Travel the tree.

The numbers are the processor numbers. This is reduction.

- n operands $\Rightarrow \log n$ steps
- How do you map?

- $n/2^i$ processors per step
- step i : if $!(id \% 2^i)$



Binomial Tree

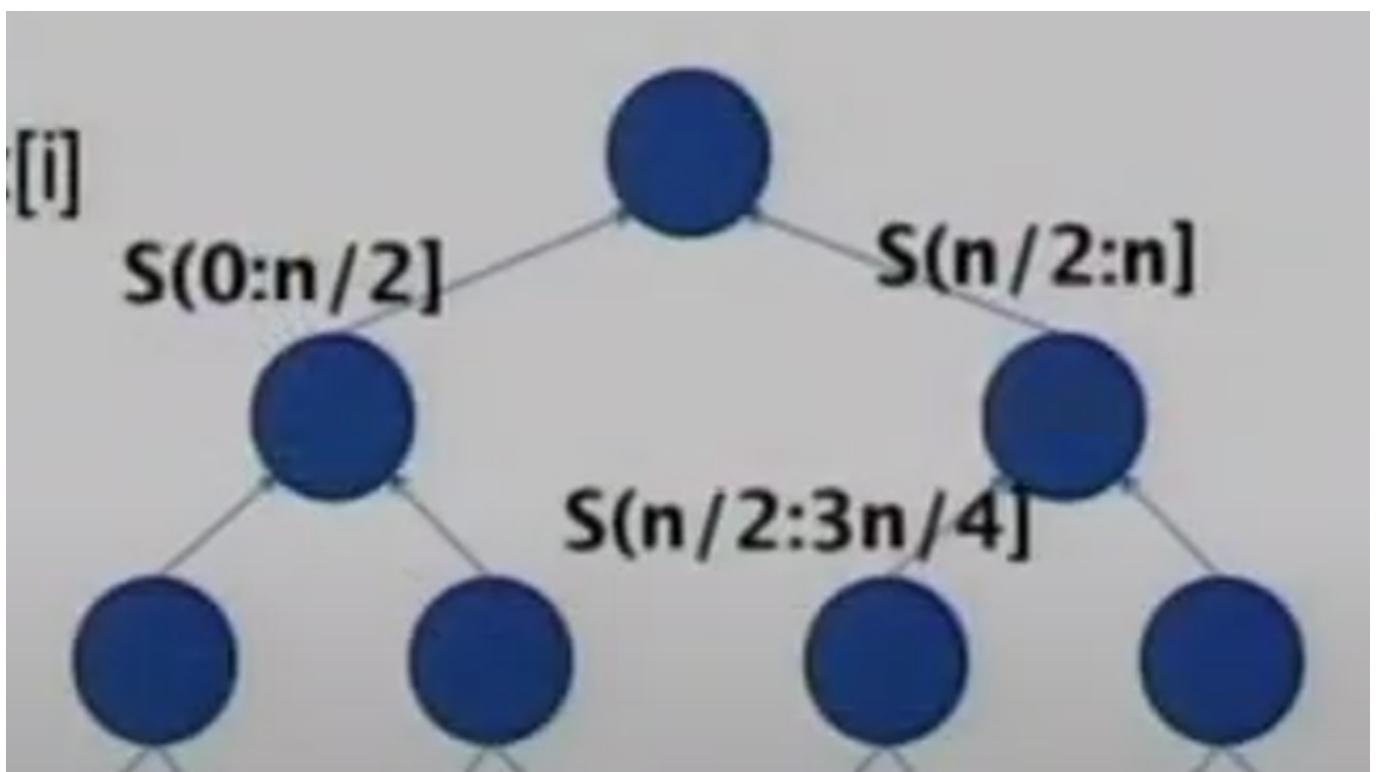
- B_0 : single node (Root)
- B_k : Root with k binomial subtrees, $B_0 \dots B_{k-1}$

The diagram illustrates a binomial tree structure across three levels. Level B_0 contains a single blue circle. Level B_1 contains two blue circles, each connected to the root of B_0 . Level B_2 contains four blue circles, each connected to one of the nodes in B_1 . Below the tree, labels B_0 , B_1 , and B_2 are positioned under their respective levels.

okay So now let us talk about somewhat non-trivial or non-obvious let us say algorithm

Each processor is dependent on itself in a binomial tree manner.

Prefix sum is sum of everything before it. $x_n = \sum_0^n x_i$. Below is prefix sum. $S(m : n]$ means sum of m th element to $n-1$ th element.



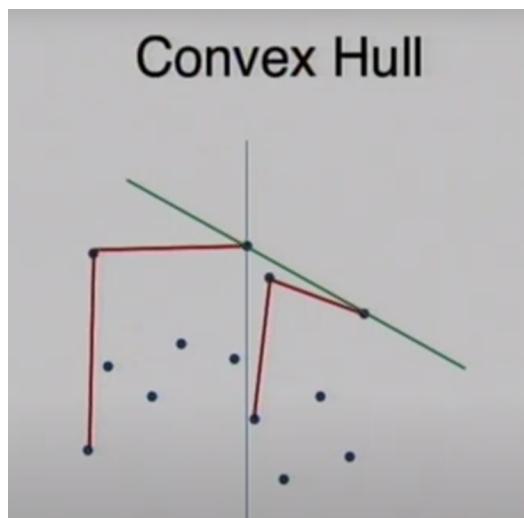
Partitioning

Merge sort: This is kinda the direct obvious implementation. Merge sort has 2 things that need to compute and combine, so

compute those together $\gamma^{(J-L)}\gamma$. HOWEVER this isn't the best parallel merge sort.

Optimal merge sort has gone over my head. Here are some key words. Split the main list into $\log \log n$ size chunks and do the merging. Complexity is $O(n * \log(\log n))$. How? ~~black magic~~. This is an example of accelerated cascading. Explained lil later. (We have taken 2 diff algos to determine that $\log \log n$ size)

Convex hull is parallelized by "splitting" the set of points. We make basic hulls for both sets of points and combine.



Accelerated cascading:

This is where we take a fast but not optimal algorithm and combine with an optimal but fast algorithm to determine the optimal chunk sizes.

Minimum number finding:

This is optimized by making a binary tree with children of nodes = $\sqrt{\# \text{nodes in tree rooted at the node's prev lvl}}$, depth = $\log \log n$. Root has 9 children, next lvl has 3 children and so on.

Recursive doubling:

Don't know.