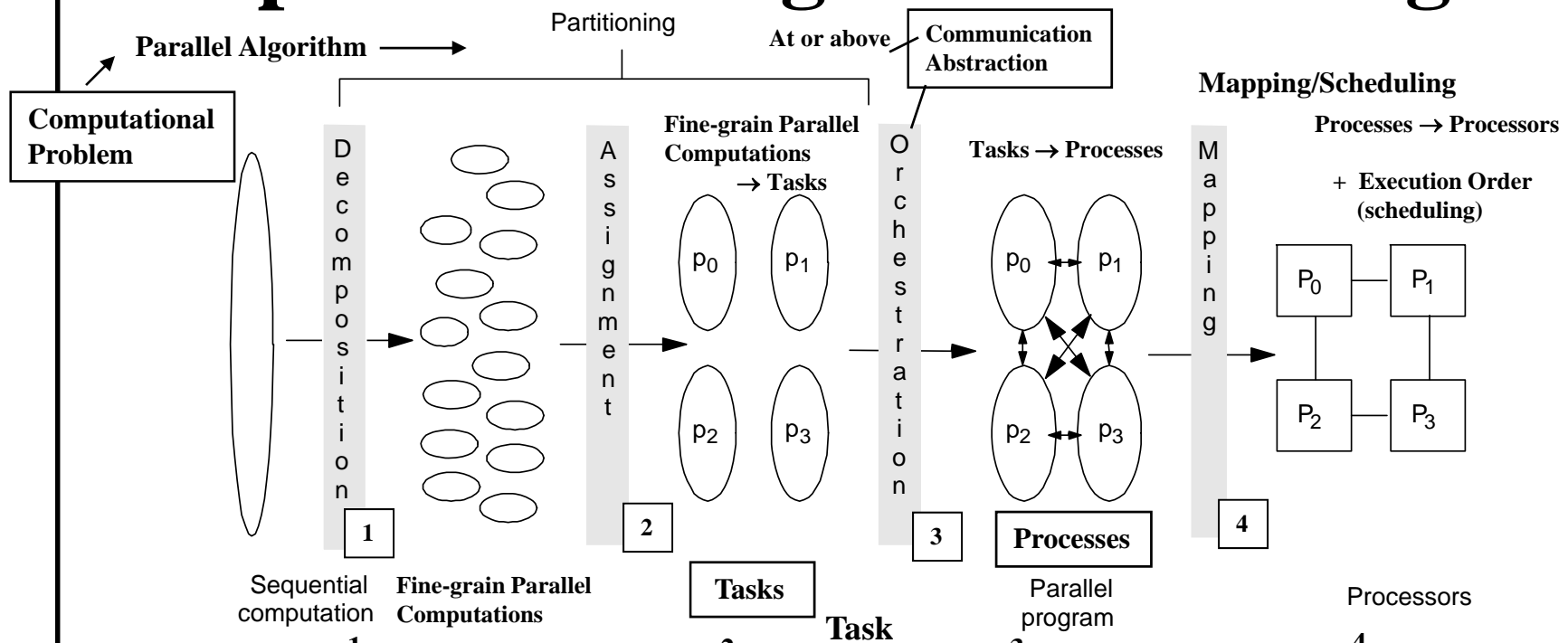


Steps in Creating a Parallel Program



- **4 steps:** ¹ **Decomposition**, ² **Assignment**, ³ **Orchestration**, ⁴ **Mapping** + Scheduling
- **Performance Goal of the steps:** Maximize parallel speedup (minimize resulting parallel) execution time by:
 - 1– Balancing computations and overheads on processors (every processor does the same amount of work + overheads).
 - 2– Minimizing communication cost and other overheads associated with each step.

Parallel Programming for Performance

A process of Successive Refinement of the steps

- Partitioning for Performance:

- Load Balancing and Synchronization Wait Time Reduction

- Identifying & Managing Concurrency

- Static Vs. Dynamic Assignment or tasking

- Determining Optimal Task Granularity

- Reducing Serialization / Synch Wait Time

- Reducing Inherent Communication

- Minimizing *communication to computation ratio* C-to-C Ratio

- Efficient Domain Decomposition Or Domain Partitioning

- Reducing Additional Overheads

- Orchestration/Mapping for Performance:

- Extended Memory-Hierarchy View of Multiprocessors For SAS

- Exploiting Spatial Locality/Reduce Artifactual Communication

- Structuring Communication

- Reducing Contention

- Overlapping Communication

Waiting time
as a result of
data dependency

or “Extra”

CMPE655 - Shaaban

Successive Refinement of Parallel Program Performance

Partitioning is possibly independent of architecture, and may be done first (initial partition):

What about number of processors?

- View machine as a collection of communicating processors
 - Balancing the workload across tasks/processes/processors.
 - Reducing the amount of inherent communication.
 - Reducing extra work to find a good assignment.
- Above three issues are conflicting.

+ Lower
C-to-C ratio

Then deal with interactions with architecture (Orchestration, Mapping) :

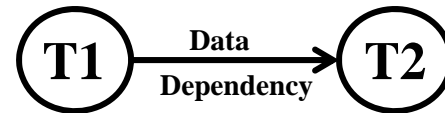
- View machine as an extended memory hierarchy:
 - Reduce artifactual (extra) communication due to architectural interactions.
 - Cost of communication depends on how it is structured (possible overlap with computation) + Hardware Architecture
- This may inspire changes in partitioning.

And algorithm?

CMPE655 - Shaaban

Partitioning for Performance

- 1 • Balancing the workload across tasks/processes:
- Reducing wait time at synchronization points needed to satisfy data dependencies among tasks.
- 2 • Reduce Overheads:
- Reducing interprocess inherent communication.
 - Reducing extra work needed to find a good assignment.



Lower
C-to-C ratio

The above goals lead to two extreme trade-offs:

- Minimize communication => run on 1 processor. One large task
 - ? \updownarrow And other parallelization overheads => extreme load imbalance.
 - Maximize load balance => random assignment of tiny tasks.
=> no control over communication.
 - A good partition may imply extra work to compute or manage it
- The goal is to compromise between the above extremes

Partitioning for Performance:

Load Balancing and “Synch Wait Time” Reduction

Limit on speedup:

Synch wait time = process/task wait time as a result of data dependency on another task
(until the dependency is satisfied)

$$Speedup_{problem}(p) \leq \frac{\text{Sequential Work}}{\underset{\text{(on any processor)}}{Max}(\text{Work on any Processor})}$$

- Work includes computing, data access and other costs (overheads).
- Not just equal work, but must be busy (computing) at same time to minimize synchronization wait time to satisfy dependencies.

Four parts to load balancing and reducing synch wait time:

1. Identify enough concurrency in decomposition.
2. Decide how to manage the concurrency (statically or dynamically).
3. Determine the granularity (task grain size) at which to exploit it.
4. Reduce serialization and cost of synchronization.

Identifying Concurrency: Decomposition

- Concurrency may be found by:

- 1 – Examining loop structure of sequential algorithm.
- 2 – Fundamental data dependencies (dependency analysis/graph).
- 3 – Exploit the understanding of the problem to devise parallel algorithms with more concurrency (e.g ocean equation solver).

i.e. 2D Grid Computation

- Software/Algorithm Parallelism Types:

1 - Data Parallelism versus 2- Functional Parallelism:

1 - Data Parallelism:

- Similar parallel operation sequences performed on elements of large data structures
 - (e.g ocean equation solver, pixel-level image processing)
- Such as resulting from parallelization of loops.
- Usually easy to load balance. (e.g ocean equation solver)
- Degree of concurrency (i.e. Max. DOP) usually increase with input or problem size. e.g $O(n^2)$ in equation solver example.

2D Grid Computation

Software/Algorithm Parallelism Types: were also covered in lecture 3 slide 33

CMPE655 - Shaaban

Identifying Concurrency (continued)

2- Functional Parallelism:

- Entire large tasks (procedures) with possibly different functionality that can be done in parallel on the same or different data.

e.g. different independent grid computations in Ocean.

- Software Pipelining: Different functions or software stages of the pipeline performed on different data:

i.e. Max.
DOP

- As in video encoding/decoding, or polygon rendering.

- Concurrency degree usually modest and does not grow with input size

- Difficult to load balance.

i.e. Problem Size

- Often used to reduce synch wait time between data parallel phases.

Most scalable parallel computations/programs:

(more concurrency as problem size increases) parallel programs:

Data parallel programs (per this loose definition)

- Functional parallelism can still be exploited to reduce synchronization wait time between data parallel phases.

Managing Concurrency: Task Assignment

Goal: Obtain an assignment with a good load balance among tasks (and processors in mapping step)

.. and Low C-to-C Ratio

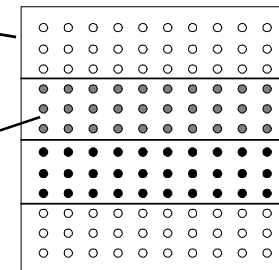
Static versus Dynamic Assignment:

Static Assignment: (e.g equation solver) i.e 2D Grid Computation

- Algorithmic assignment usually based on input data ; does not change at run time.
- Low run time overhead.
- Computation must be predictable.
- Preferable when applicable (lower overheads).

of computations into tasks at compilation time

e.g. strip of n/p rows per task



Example 2D Ocean Equation Solver

At Compilation Time

Dynamic Task Assignment: Or dynamic tasking

- Needed when computation not fully predictable.
- Adapt partitioning at run time to balance load on processors.
- Can increase communication cost and reduce data locality.
- Can increase run time task management overheads. Counts as extra work

At Run Time

CMPE655 - Shaaban

Dynamic Task Assignment/Mapping

Profile-based (semi-static):

Initial partition

- Profile (algorithm) work distribution initially at runtime, and repartition dynamically.
- Applicable in many computations, e.g. Barnes-Hut, → (simulating galaxy evolution) some graphics.

Dynamic Tasking:

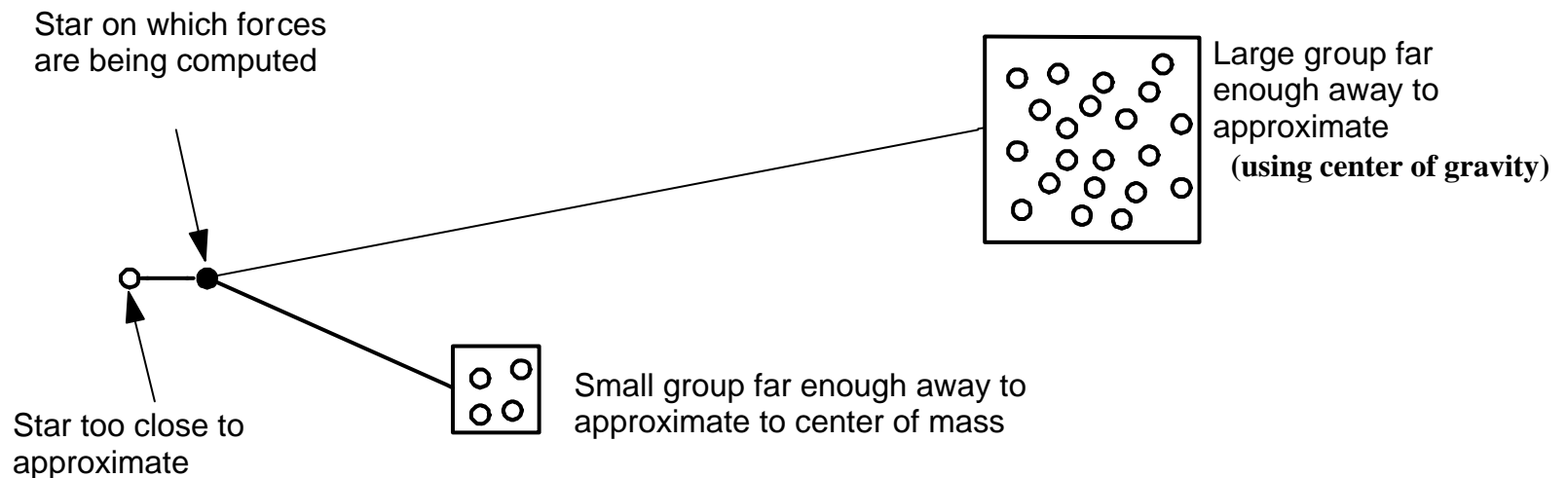
- Deal with unpredictability in parallel computation, program or environment (e.g. Ray tracing).
 - Computation, communication, and memory system interactions
 - Multiprogramming and heterogeneity of processors
 - Used by runtime systems and OS too.
- Pool (queue) of tasks: Processors take and add tasks to pool until parallel computation is done.
- e.g. “self-scheduling” of loop iterations (shared loop counter).

Simulating Galaxy Evolution

(Gravitational N-Body Problem)

- Simulate the interactions of many stars evolving over time
- Computing forces is expensive
 - $O(n^2)$ brute force approach
- Hierarchical Methods (e.g. Barnes-Hut) take advantage of force law: G (*center of mass*)

$$\frac{m_1 m_2}{r^2}$$

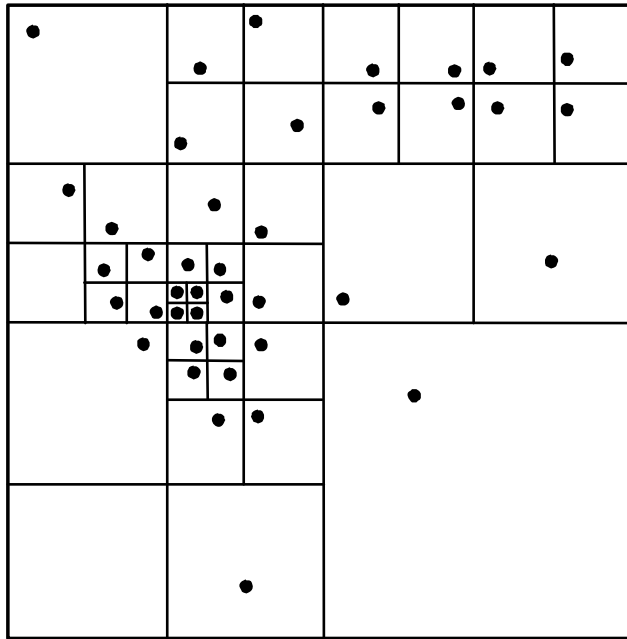


- Many time-steps, plenty of concurrency across stars within one

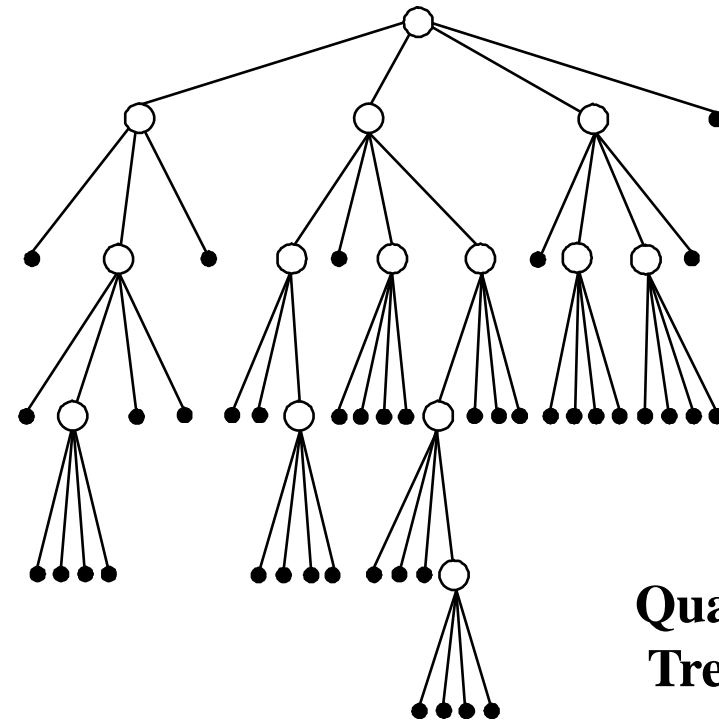
Gravitational N-Body Problem: Barnes-Hut Algorithm

- To parallelize problem: Groups of bodies partitioned among processors. Forces communicated by messages between processors.
 - Large number of messages, $O(N^2)$ for one iteration.
- Solution: Approximate a cluster of distant bodies as one body with their total mass
- This clustering process can be applied recursively.
- Barnes_Hut: Uses divide-and-conquer clustering. For 3 dimensions:
 - Initially, one cube contains all bodies
 - Divide into 8 sub-cubes. (4 parts in two dimensional case).
 - If a sub-cube has no bodies, delete it from further consideration.
 - If a cube contains more than one body, recursively divide until each cube has one body
 - This creates an oct-tree which is very unbalanced in general.
 - After the tree has been constructed, the total mass and center of gravity is stored in each cube.
 - The force on each body is found by traversing the tree starting at the root stopping at a node when clustering can be used.
 - The criterion when to invoke clustering in a cube of size $d \times d \times d$:
$$r \geq d/\theta$$
 - r = distance to the center of mass
 - θ = a constant, 1.0 or less, opening angle
 - Once the new positions and velocities of all bodies is computed, the process is repeated for each time period requiring the oct-tree to be reconstructed (repartition dynamically)

Two-Dimensional Barnes-Hut



(a) The spatial domain



(b) Quadtree representation

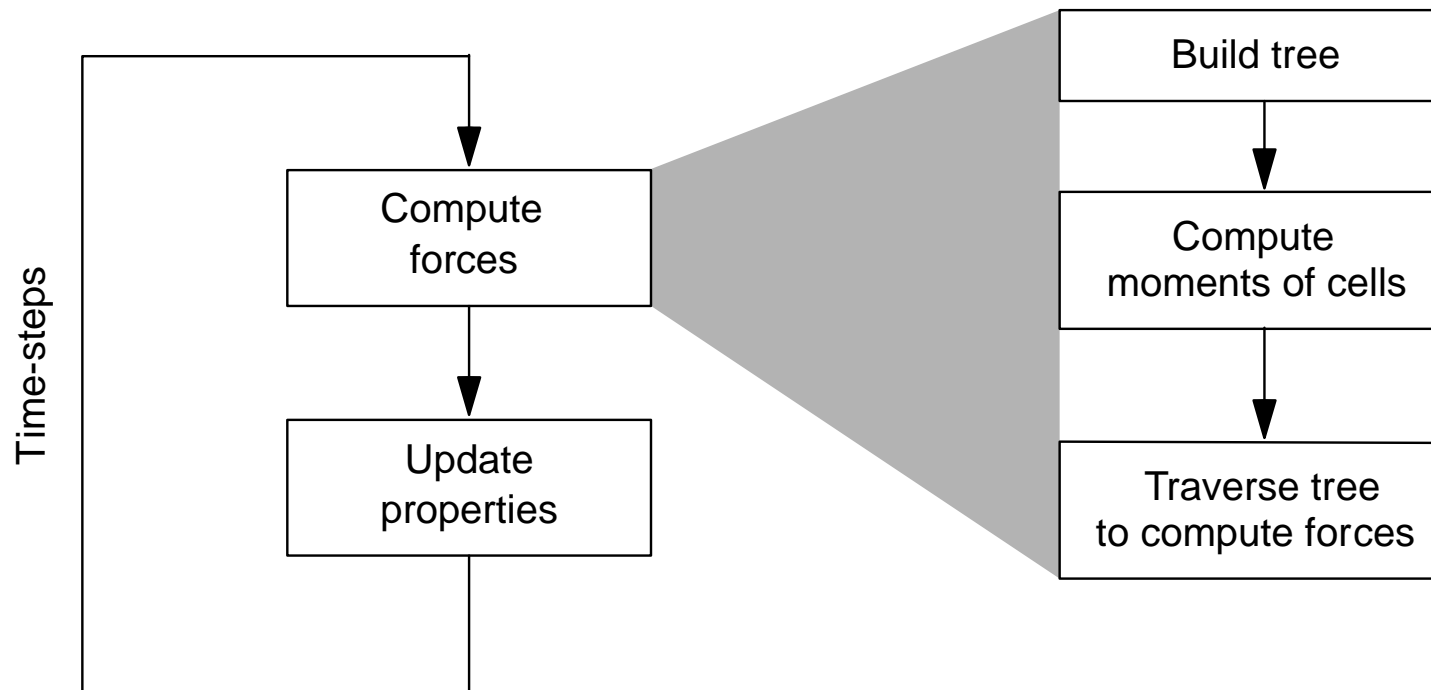
Quad
Tree

Recursive Division of Two-dimensional Space

Locality Goal:

Bodies close together in space should be on same processor

Barnes-Hut Algorithm



- **Main data structures: array of bodies, of cells, and of pointers to them**
 - **Each body/cell has several fields: mass, position, pointers to others**
 - **pointers are assigned to processes**

The Need For Dynamic Tasking: Rendering Scenes by Ray Tracing

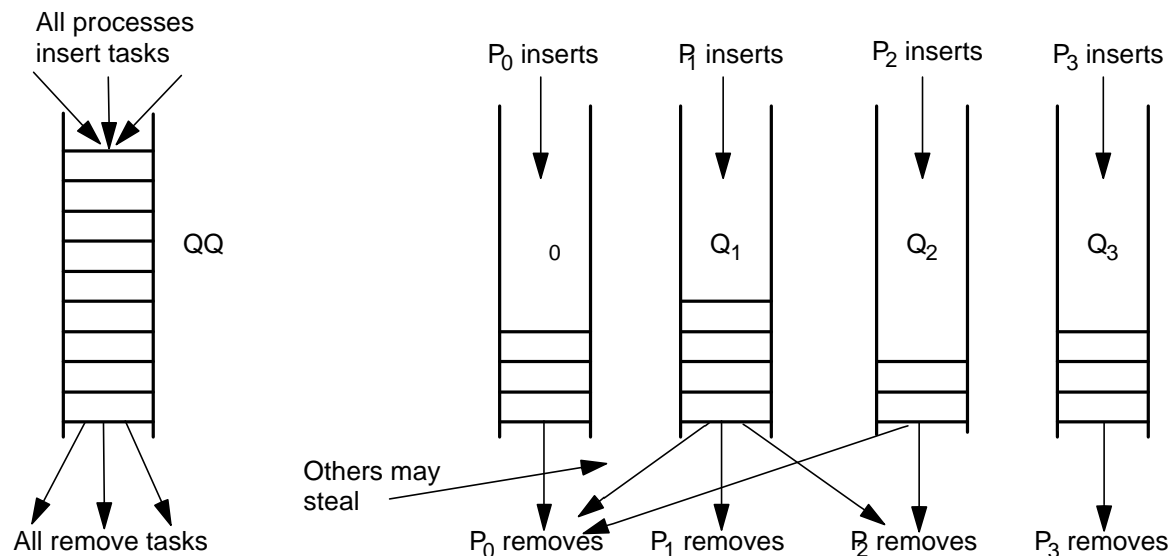
- Shoot rays into a scene through pixels in image plane.
- Follow their paths: To/from light sources/reflected light
 - They bounce around as they strike objects:
 - They generate new rays:
 - Resulting in a ray tree per input ray and thus more computations (tasks).
- Result is color and opacity for that pixel.
- Parallelism across rays.
 - Parallelism here is unpredictable statically.
 - Dynamic tasking needed for load balancing.

Dynamic Tasking with Task Queues

Centralized versus distributed queues.

Task stealing with distributed queues.

- Can compromise communication and data locality (e.g in SAS), and increase synchronization wait time.
- Whom to steal from, how many tasks to steal, ...
- Termination detection (all queues empty).
- Load imbalance possible related to size of task.
 - Many small tasks usually lead to better load balance

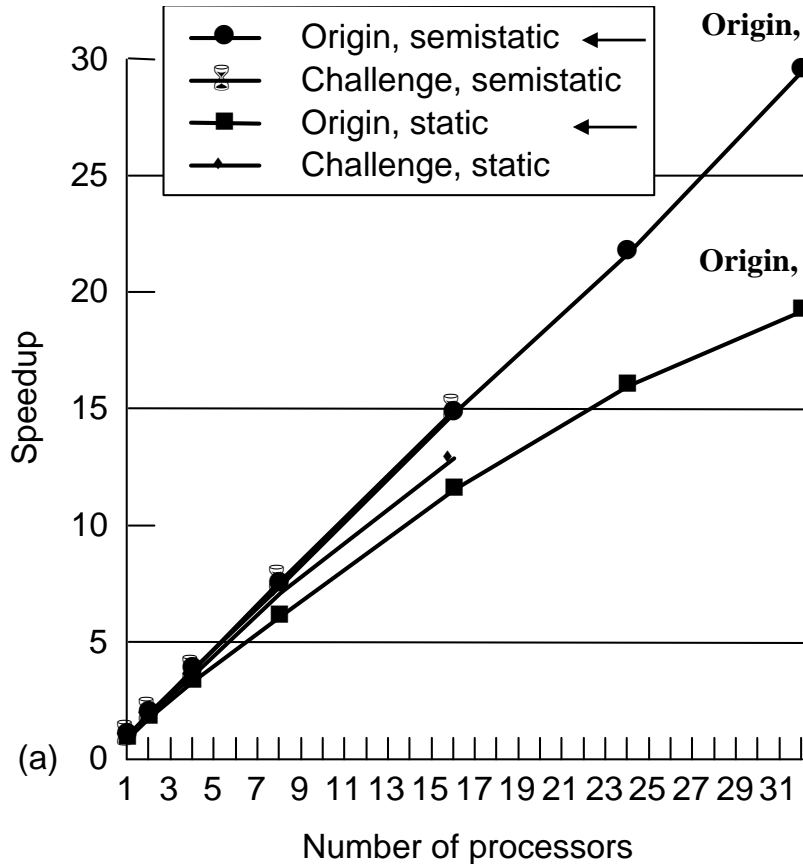


Centralized (One) Task Queue

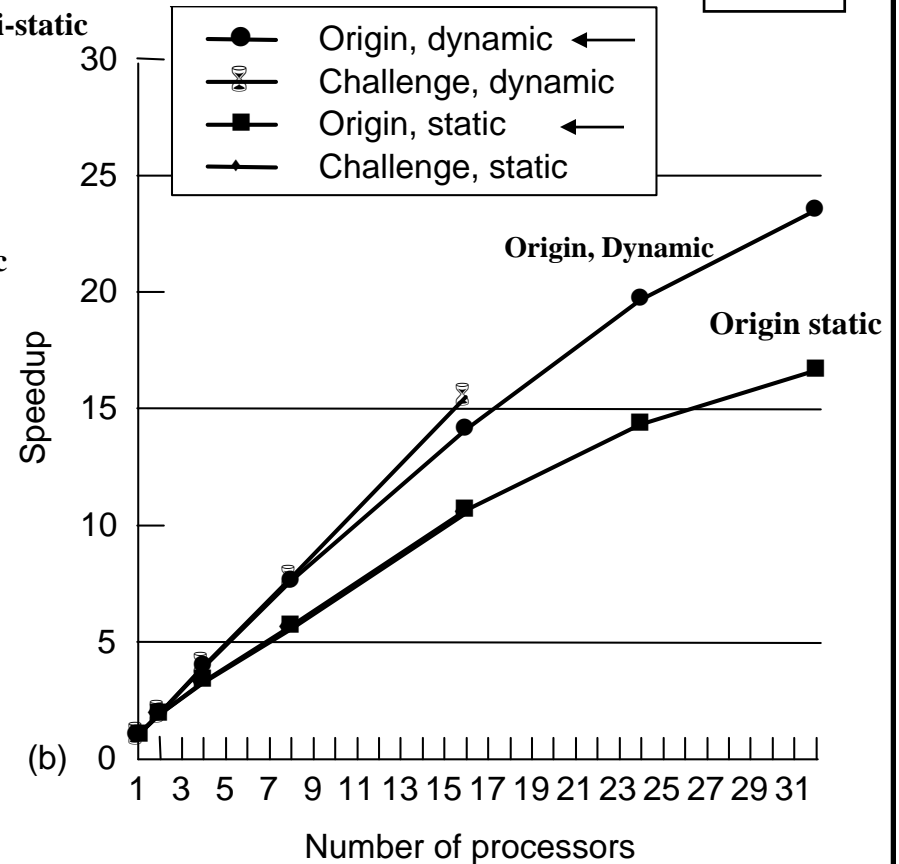
Distributed Task Queues (one per process)

Performance Impact of Dynamic Assignment

On SGI Origin 2000 (cache-coherent shared distributed memory):



Barnes-Hut 512k particle
(N-Body Problem)



Ray tracing

NUMA

Partitioning for Performance:

Task

Assignment: Determining Task Size/Granularity

Recall that parallel task granularity:

Amount of work or computation associated with a task.

General rule:

Always ?

- Coarse-grained => Often less load balance
Larger/fewer tasks less communication and other overheads
- Fine-grained => more overhead; often more communication, contention
Smaller/more tasks But potentially better load balance

Communication, contention actually more affected by mapping to processors, not just task size only.

- Other overheads are also affected by task size too, particularly with dynamic mapping (tasking) using task queues:
 - Small tasks -> More Tasks -> More dynamic mapping overheads.

Or task queue work units

A task only executes on one processor to which it has been mapped or allocated

CMPE655 - Shaaban

Reducing Serialization/Synch Wait Time

Requires careful assignment and orchestration (and scheduling ?)

Reducing Serialization/Synch wait time in Event synchronization:

- Reduce use of conservative synchronization e.g. : i.e Ordering
 - Fine point-to-point synchronization instead of barriers (if possible),
 - or reduce granularity of point-to-point synchronization (specific elements instead of entire data structure).
- But fine-grained synch more difficult to program, more synch operations.

Reducing Serialization in Mutual exclusion: Lock Enter Critical Section Exit Unlock

1- Separate locks for separate data

- e.g. locking records in a database instead of locking entire database: lock per process, record, or field
- Lock per task in task queue, not per queue
- Finer grain => less contention/serialization, more space, less reuse

2- Smaller, less frequent critical sections

e.g use of local difference in 2D Grid example

- No reading/testing in critical section, only modification
- e.g. searching for task to dequeue in task queue, building tree etc.

3- Stagger critical sections in time (on different processors).

i.e critical section entry occur at different times

Partitioning for Performance:

Implications of Load Balancing/Synch Time Reduction

Extends speedup limit expression to:

$$Speedup_{problem}(p) \leq \frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time)}}_{(on\ any\ processor)}$$

Generally load balancing is the responsibility of software

For dynamic tasking

But: Architecture can support task stealing and synch efficiently:

- *Fine-grained* communication, *low-overhead access* to queues
 - Efficient support allows smaller tasks, better load balancing
 - *Naming* logically shared data in the presence of task stealing
 - Need to access data of stolen tasks, esp. multiple-stolen tasks
- => Hardware shared address space advantageous here
- Efficient support for point-to-point communication.
 - Software layers + hardware (network) support.
+ CA

Partitioning for Performance:

Reducing Inherent Communication

Inherent Communication : communication between tasks inherent in the problem/parallel algorithm for a given partitioning/assignment (to tasks)

Measure: *communication to computation ratio*
(*c-to-c ratio*)

Focus here is on reducing interprocess communication inherent in the problem:
i.e inherent communication

- Determined by assignment of parallel computations to tasks/processes.
- Minimize c-to-c ratio while maintaining a good load balance among tasks/processes.
- Actual communication can be greater than inherent communication.

e.g.

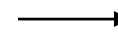
As much as possible, assign tasks that access same data to same process (and processor later in mapping).

Processor Affinity

Important in SAS NUMA Architectures

- Optimal solution (partition) to reduce communication and achieve an optimal load balance is NP-hard in the general case.
- Simple heuristic partitioning solutions may work well in practice:
 - Due to specific dependency structure of applications.
 - Example: Domain decomposition

Next



or domain partitioning

CMPE655 - Shaaban

Example Assignment/Partitioning Heuristic:

Domain: Physical domain of problem or input data set

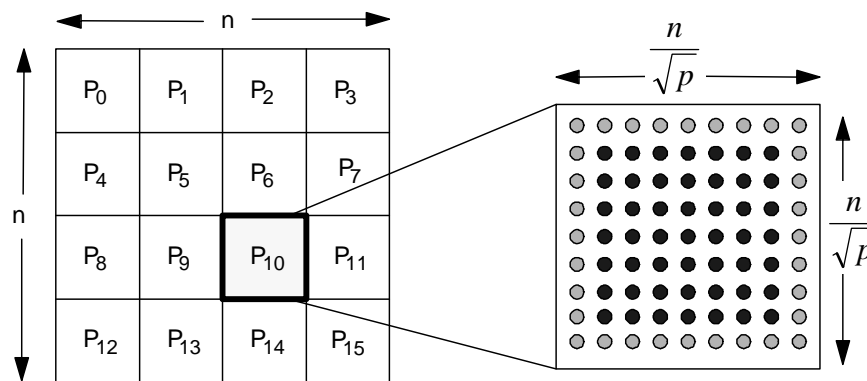
Domain Decomposition

- Initially used in data parallel scientific computations such as (Ocean) and pixel-based image processing to obtain a good load balance and c-to-c ratio. *and other usually predictable computations tied to a physical domain/data set*

How?

The task assignment is achieved by decomposing the physical domain or data set of the problem. *Such assignment often done statically for predictable computations*

- Exploits the local-biased nature of physical problems
 - Information requirements often short-range
 - Or long-range but fall off with distance
- Simple example: Nearest-neighbor 2D grid computation (as in ocean example)



Block Decomposition Or assignment

$$Communication = \frac{4n}{\sqrt{p}} \quad Computation = \frac{n^2}{p}$$

$$C-to-C = \frac{4 \times \sqrt{p}}{n}$$

comm-to-comp ratio = Perimeter to Area (area to volume in 3-d)

- Depends on n, p : decreases with n , increases with p

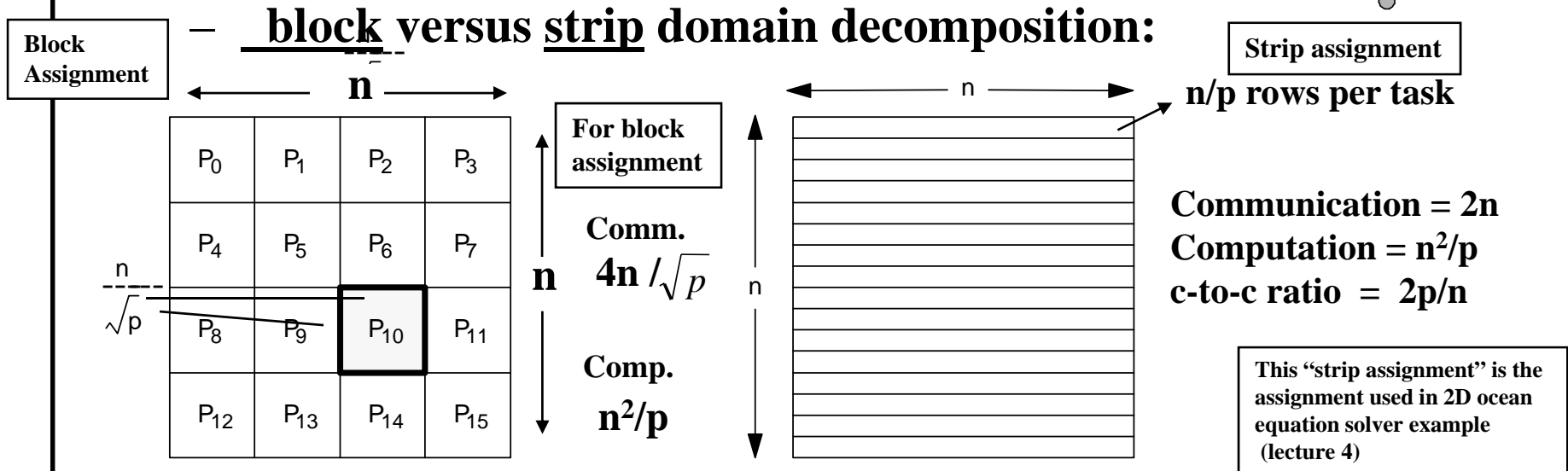
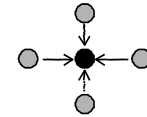
CMPE655 - Shaaban

$n \times n$ Grid p = Number of tasks/processes here = $p = 4 \times 4 = 16$

Domain Decomposition (continued)

Best domain decomposition depends on information requirements

Nearest neighbor example: i.e group or "strip" of (contiguous) rows



Block Decomposition

Strip (Group of rows) Decomposition

Comm-to-comp ratio: $\frac{4 \times \sqrt{p}}{n}$ for block, $\frac{2 \times p}{n}$ for strip

Which C-to-C ratio is better?

Application dependent: strip may be better in some cases

n x n Grid p processors Often n >> p

CMPE655 - Shaaban

Finding a Domain Decomposition

Four possible methods:

More
Work
↓

1 • Static, by inspection:

- Computation must be predictable: e.g grid example above, and Ocean

Not input data dependent

and low-level (pixel-based) image processing

2 • Static, but not by inspection:

- Input-dependent, require analyzing input structure
 - Before start of computation once input data is known.
- E.g sparse matrix computations, data mining

Characterized by non-uniform data/computation distribution

3 • Semi-static (periodic repartitioning):

- Characteristics change but slowly; e.g. Barnes-Hut

N-Body
Problem

4 • Static or semi-static, with dynamic task stealing

- Initial decomposition based on domain, but highly unpredictable computation; e.g ray tracing

May require the utilization of one or more task queues to obtain a good load balance

CMPE655 - Shaaban

Implications of Communication

- Architects must examine application latency/bandwidth needs
- If denominator in c-to-c is computation execution time, ratio gives average BW needs per task.
- If denominator in c-to-c is operation count, gives extremes in impact of latency and bandwidth
 - Latency: assume no latency hiding.
 - Bandwidth: assume all latency hidden.
 - Reality is somewhere in between.
- Actual impact of communication depends on structure and cost as well:

Communication Cost = Time added to parallel execution time as a result of communication

From lecture 2

$$\text{Speedup} \leq \frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time + Comm Cost)}}_{\text{(on any processor)}}$$

→ Need to keep communication balanced across processors as well.

c-to-c = communication to computation ratio

CMPE655 - Shaaban

Partitioning for Performance: Reducing Extra Work (Overheads)

Must also be balanced
among all processors

- **Common sources of extra work (mainly orchestration):**

- Computing a good partition (at run time):
e.g. partitioning in Barnes-Hut or sparse matrix
- Using redundant computation to avoid communication.
- Task, data distribution and process management overhead
 - Applications, languages, runtime systems, OS
- Imposing structure on communication:
 - Coalescing (combining) messages, allowing effective naming

- **Architectural Implications:**

More on this a bit later in the lecture

- Reduce by making communication and orchestration efficient
(e.g hardware support of primitives ?)

$$\text{Speedup} \leq \frac{\text{Sequential Work}}{\text{Max (Work + Synch Wait Time + Comm Cost + Extra Work)}} \\ \text{(on any processor)}$$

Summary of Parallel Algorithms Analysis

- Requires characterization of multiprocessor system and algorithm requirements.
- Historical focus on algorithmic aspects: partitioning, mapping
- In PRAM model: data access and communication are *free*
 - Only load balance (including serialization) and extra work matter

For PRAM:

$$\text{Speedup}_{\text{PRAM}} \leq \frac{\text{Sequential Instructions}}{\text{Max (Instructions + Synch Wait Time + Extra Instructions)}}$$

(on any processor)

Or Work

extra work/computation not in sequential version

PRAM
Advantages/
Disadvantages

- Useful for parallel algorithm development, but possibly unrealistic for real parallel program performance evaluation.
 - Ignores communication and also the imbalances it causes
 - Can lead to poor choice of partitions as well as orchestration when targeting real parallel systems.

Limitations of Parallel Algorithm Analysis

i.e communication between tasks inherent in the problem/parallel algorithm for a given partitioning/assignment (to tasks)

Possible metric for inherent communication: C-to-C Ratio

- **Inherent communication in a parallel algorithm is not the only communication present:**

- **Artifactual “extra” communication caused by program implementation and architectural interactions can even dominate.**

- Thus, actual amount of communication may not be dealt with adequately *i.e If artifactual communication is not accounted for*

- **Cost of communication determined not only by amount:** +

- Also how communication is structured and overlapped. +

- Cost of communication (primitives) in system +

- Software related and hardware related (network) — including CA

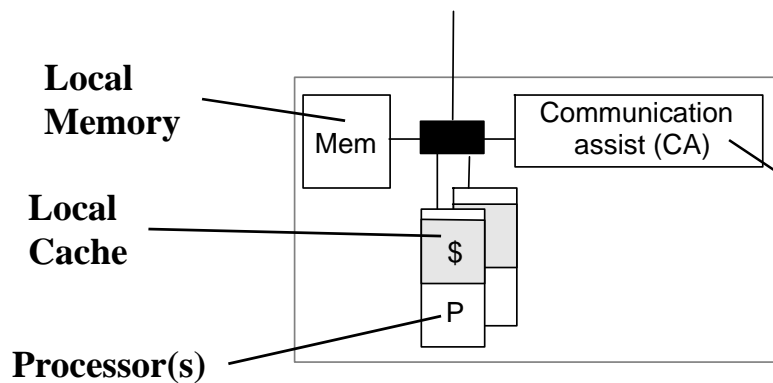
- Both are architecture-dependent, and addressed in orchestration step.

Generic Multiprocessor Architecture

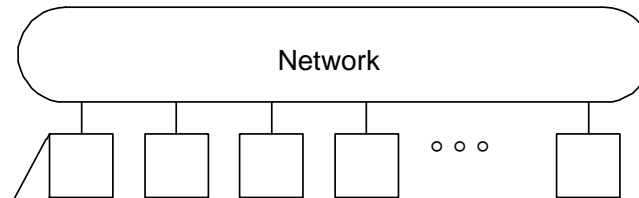
If SAS is natively supported by this generic architecture:

→ NUMA

(Distributed Shared memory Architecture)



Scalable network.



Nodes

CA may support SAS in hardware or just message-passing

Computing Nodes:

processor(s), memory system, plus *communication assist (CA)*:

- Network interface and communication controller.

Scalable Network.

Extended Memory-Hierarchy View of

For NUMA SAS
Parallel Architectures

Generic Multiprocessors

SAS support
Assumed



- **Levels in extended hierarchy:**

- ¹ Registers, ² caches, ³ local memory, ⁴ remote memory (over network) and caches

- **Glued together by communication architecture**

i.e Minimum size of
data transferred
between levels

- Levels communicate at a certain granularity of data transfer. (e.g. Cache blocks, pages etc.)

extended

- **Need to exploit spatial and temporal locality in hierarchy**

- Otherwise artifactual (extra) communication may also be caused

Why?

- Especially important since **communication is expensive**

Over network

This extended hierarchy view is more useful in distributed shared memory (NUMA) parallel architectures

CMPE655 - Shaaban

Extended Hierarchy

- **Idealized view:** local cache hierarchy + single main memory
- **But reality is more complex:**
 - **Centralized Memory:** + caches of other processors
 - **Distributed Memory:** some local, some remote; + network topology + local and remote caches
 - **Management of levels:**
 - Caches managed by hardware
 - Main memory depends on programming model:
 - SAS: data movement between local and remote transparent
 - Message passing: explicit by sending/receiving messages.
 - **Improve performance through architecture or program locality (maximize local data access).**

Otherwise artifactual “extra” communication is created

This extended hierarchy view is more useful in distributed shared memory parallel architectures

CMPE655 - Shaaban

Artifactual Communication in Extended Hierarchy

Accesses not satisfied in local portion cause communication

– Inherent Communication, implicit or explicit, causes transfers:

- Determined by parallel algorithm/program partitioning

C-to-C
Ratio

– Artifactual “Extra” Communication:

- Determined by program implementation and architecture interactions
- Poor allocation of data across distributed memories: data accessed heavily used by one node is located in another node’s local memory.
- Unnecessary data in a transfer: More data communicated in a message than needed.
- Unnecessary transfers due to system granularities (cache block size, page size).
- Redundant communication of data: data value may change often but only last value needed.
- Finite replication capacity (in cache or main memory)

For
NUMA

Why?

Causes of
Artifactual
“extra”
Communication

i.e zero or no extra
communication

Need to load data again over network

For replication

– Inherent communication assumes ¹unlimited capacity, ²small transfers, ³perfect knowledge of what is needed.

– More on artifactual communication later; first consider replication-induced further

CMPE655 - Shaaban

As defined earlier: Inherent Communication : communication between tasks inherent in the problem/parallel algorithm for a given partitioning/ assignment (to tasks)

Extra Communication and Replication

replication i.e. cache is filled to capacity

- **Extra Comm. induced by finite capacity is most fundamental artifact:**
 - Similar to cache size and miss rate or memory traffic in uniprocessors.
 - Extended memory hierarchy view useful for this relationship
- **View as three level hierarchy for simplicity**
 - Local cache, local memory, remote memory (ignore network topology).
- **Classify “misses” in “cache” at any level as for uniprocessors**
 - 1 • *Compulsory* or *cold* misses (no size effect)
 - 2 • *Capacity* misses (yes)
 - 3 • *Conflict* or *collision* misses (yes)
 - 4 • *Communication* or *coherence* misses (yes) 4 Cs
 - Each may be helped/hurt by large transfer granularity (spatial locality). e.g. Cache Block Size or Page Size

New C

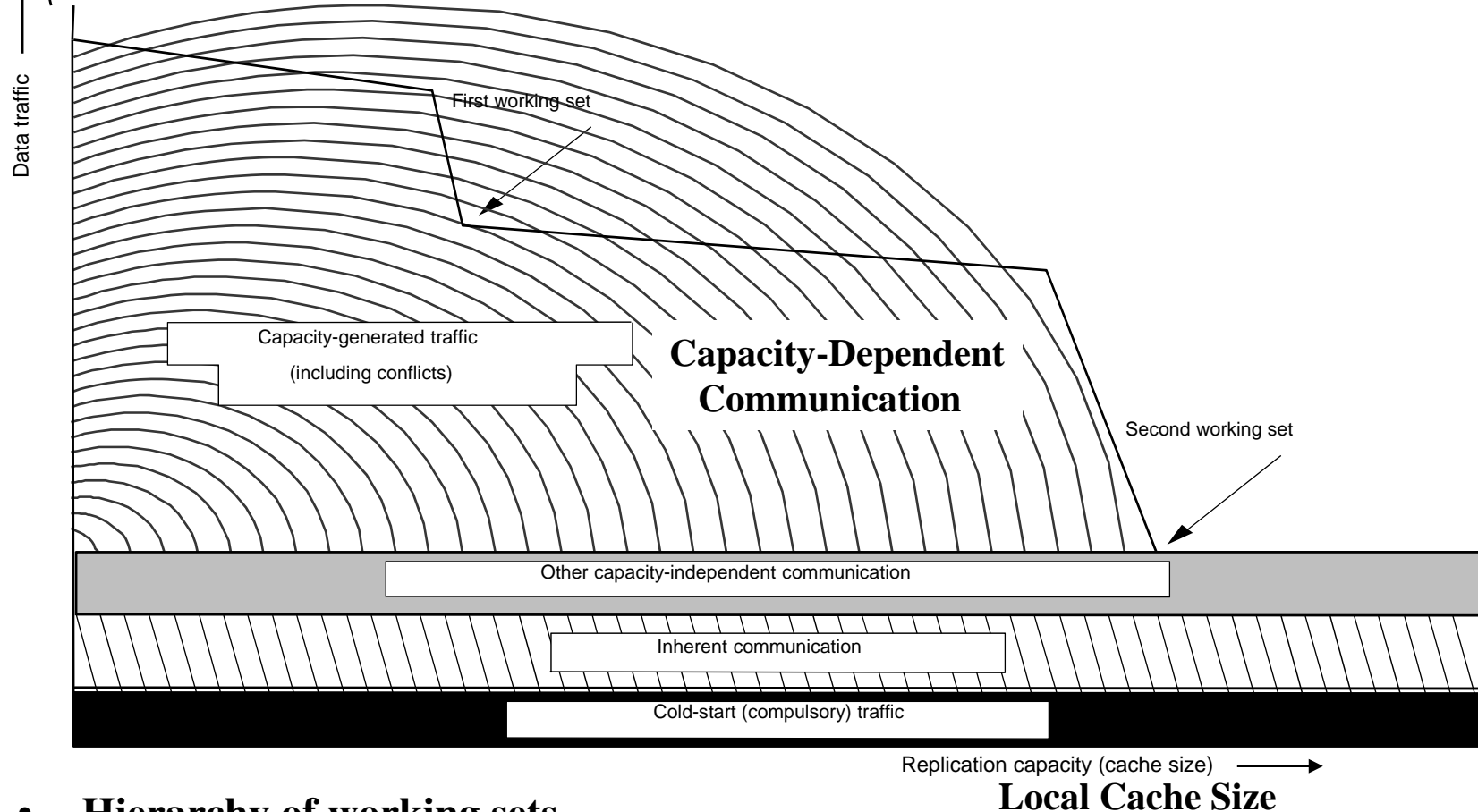
i.e misses that result in extra communication over the network

Distributed shared memory (NUMA) parallel architecture implied here

CMPE655 - Shaaban

Working Set Perspective

The data traffic between a cache and the rest of the system and components data traffic as a function of local cache size



- Hierarchy of working sets
- Traffic from any type of miss can be local or non-local (communication)

Distributed shared memory/SAS parallel architecture assumed here

CMPE655 - Shaaban

Orchestration for Performance

- **Reducing amount of communication:**
 - **Inherent:** change logical data sharing patterns in algorithm
 - Reduce c-to-c-ratio. Go back and change task assignment/partition
 - **Artifactual:** exploit spatial, temporal locality in extended hierarchy. For SAS NUMA machines
 - Techniques often similar to those on uniprocessors
- **Structuring communication to reduce cost:**
 - e.g overlap communication with computation or other communication
- We'll examine techniques for both...

Orchestration
Related

Reducing Artifactual Communication

- Message Passing Model:

- Communication and replication are both explicit.
- Even artifactual communication is in explicit messages
 - e.g. more data sent in a message than actually needed

As seen earlier

- Shared Address Space (SAS) Model:

- More interesting from an architectural perspective
- Occurs transparently due to interactions of program and system:
 - Caused by sizes of allocation and granularities in extended memory hierarchy (e.g. Cache block size, page size).

i.e. Artifactual Comm.

+ poor data allocation (NUMA)

- Next, we use shared address space to illustrate issues
→ (distributed memory SAS - NUMA)

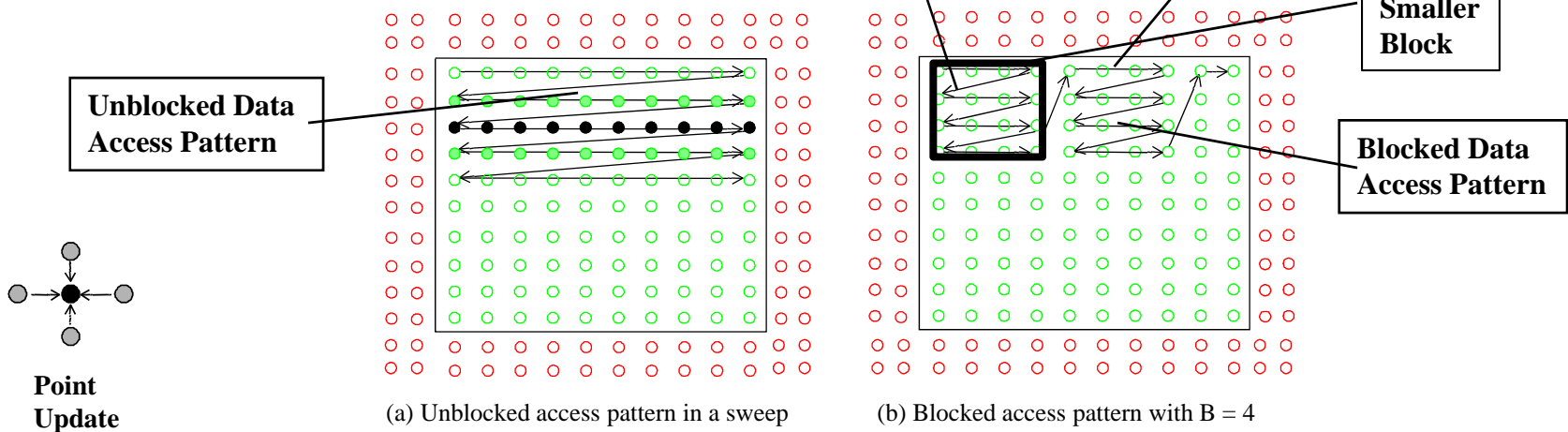
Reducing Artifactual “Extra” Communication

Exploiting Temporal Locality

- Structure algorithm so working sets map well to hierarchy
 - Often techniques to reduce inherent communication do well here
 - **Schedule tasks for data reuse once assigned** *To increase temporal locality*
- Multiple data structures in same phase
 - e.g. database records: local versus remote

i.e. 2D Grid Computation

- **Solver example: blocking** (or blocked data access pattern)



- More useful when $O(n^{k+1})$ computation on $O(n^k)$ data *i.e. computation with data reuse*
 - Many linear algebra computations (factorization, matrix multiply)

Blocked assignment assumed here

CMPE655 - Shaaban

Reducing Artifactual “Extra” Communication

Exploiting Spatial Locality

- Besides capacity, granularities are important:

- Granularity of allocation (e.g. page size)
- Granularity of communication or data transfer
- Granularity of coherence (e.g. cache block size)

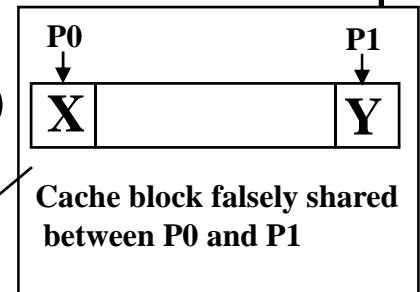
Larger “granularity”
when farther from
processor

- Major spatial-related causes of artifactual communication:

- Conflict misses
- Data distribution/layout (allocation granularity)
- Fragmentation (communication granularity)
- False sharing of data (coherence granularity)

e.g. page size

More data
sent than
needed due
to minimum
message
size



- All depend on how spatial access patterns interact with data structures/architecture:

Fix?

- Fix problems by modifying data structures, or layout/alignment (as shown in example next)

- Examine later in context of architectures

- One simple example here: data distribution in SAS solver

i.e. 2D Grid
Computation

Next

CMPE655 - Shaaban

Distributed memory (NUMA) SAS assumed here

**Reducing Artifactual
“Extra” Communication**

Spatial Locality Example

- Repeated sweeps over elements of 2D grid, block assignment, Shared address space;
- In Distributed memory: A memory page is allocated in one nodes memory (NUMA SAS)
- Natural 2D versus higher-dimensional (4D here) array representation

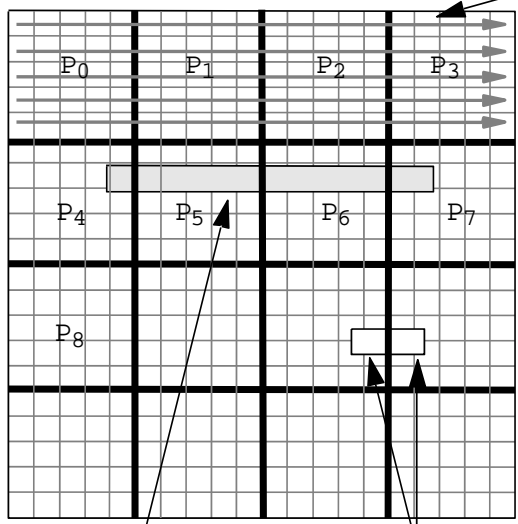
i.e granularity of data allocation

Ex: (1024, 1024)

Contiguity in memory layout

Ex: (4, 4, 256, 256)

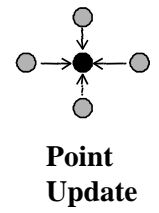
2D Array Representation



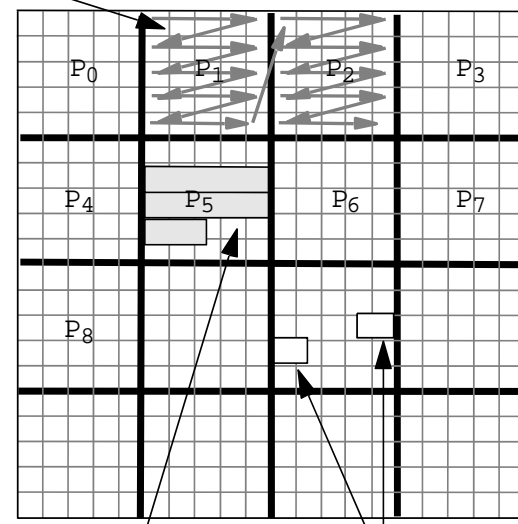
Block Assignment Used Here (for both)

Page straddles partition boundaries: difficult to distribute memory well

Cache block straddles partition boundary



4D Array Representation



Page does not straddle partition boundary

Cache block is within a partition

Two-Dimensional (2D) Array

Four-Dimensional (4D) Array

(Generates more artifactual “extra” communication)

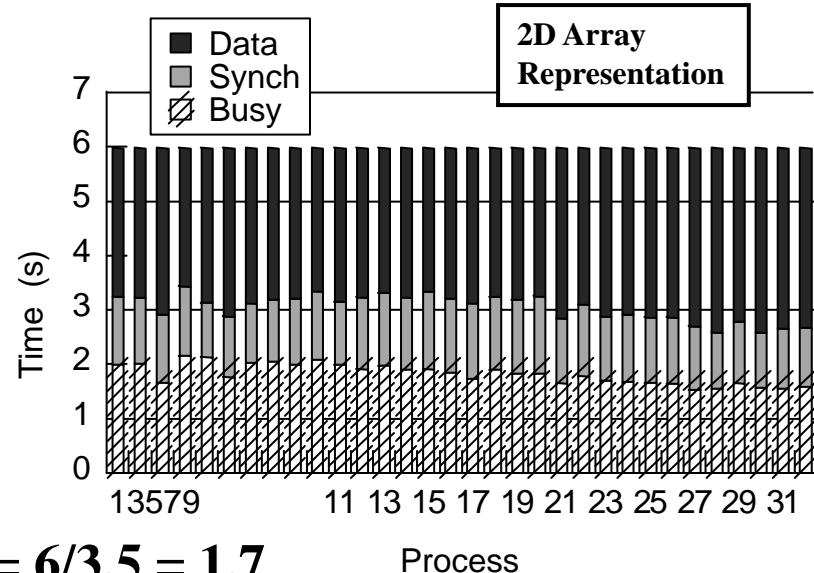
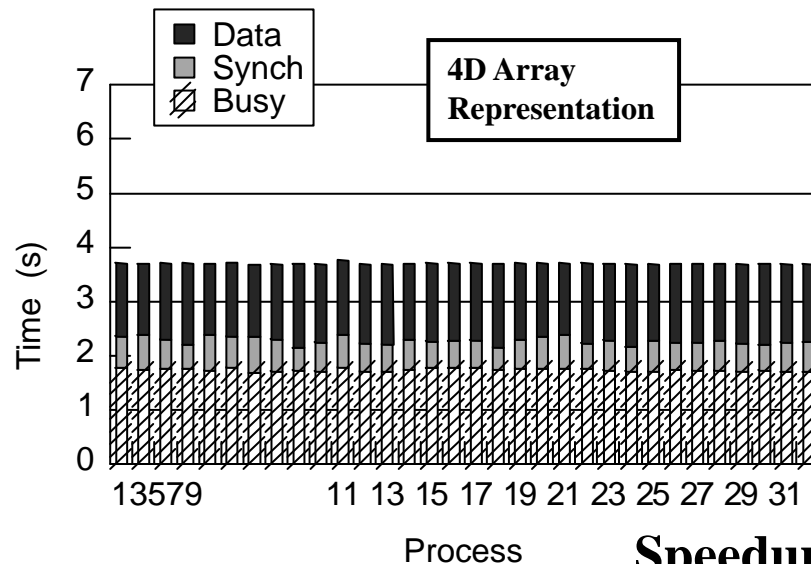
SAS assumed here

Performance Comparison Next

CMPE655 - Shaaban

Execution Time Breakdown for Ocean on a 32-processor Origin2000

1026 x 1026 grids with block partitioning on 32-processor Origin2000



$$\text{Speedup} = 6/3.5 = 1.7$$

Four-dimensional (4D) arrays

Two-dimensional (2D) arrays

- 4D grids much better than 2D, despite very large caches on machine (4MB L2 cache)

Thus less replication capacity

- data distribution is much more crucial on machines with smaller caches

- Major bottleneck in this configuration is time waiting at barriers

- imbalance in memory stall times as well

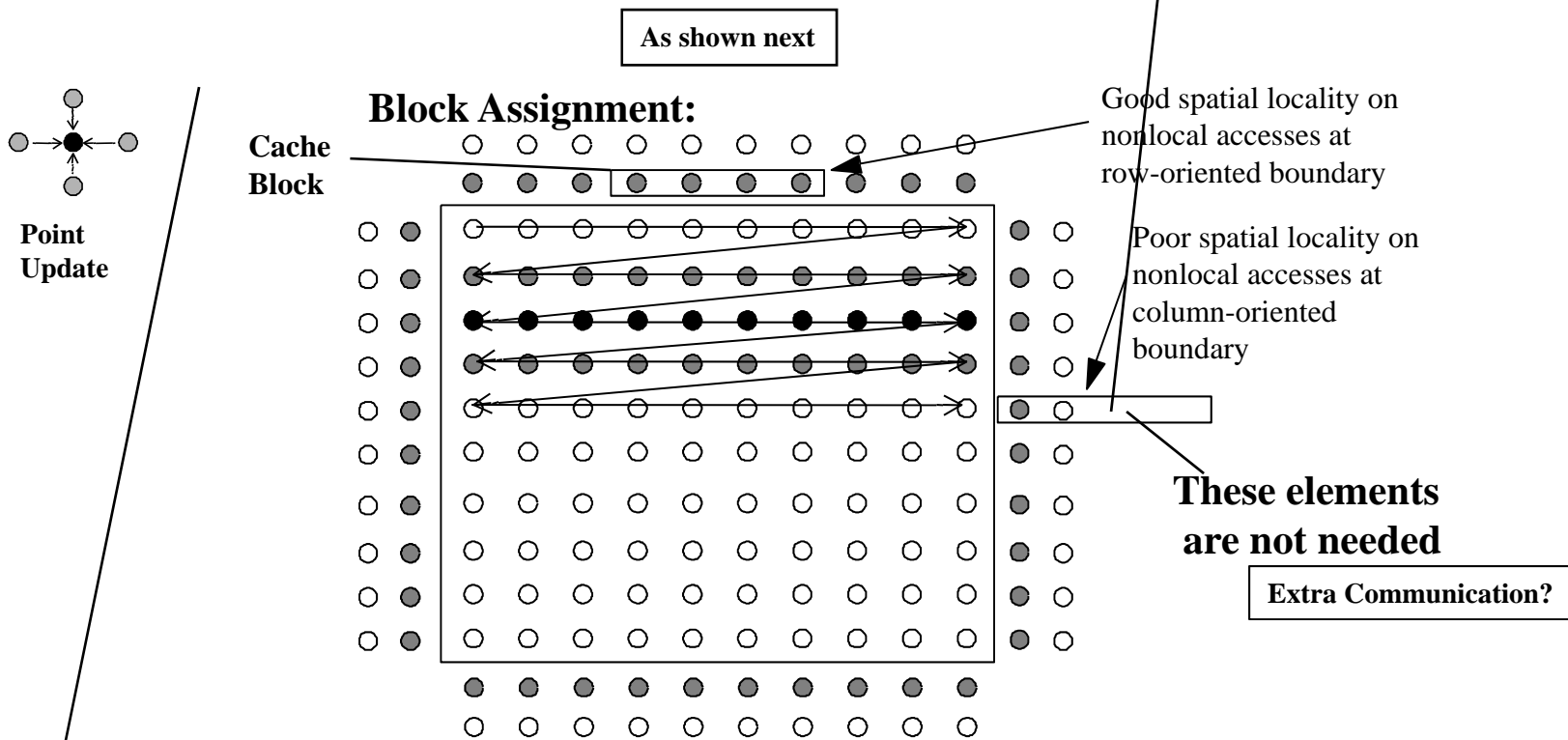
CMPE655 - Shaaban

Tradeoffs with Inherent Communication

i.e block assignment

Partitioning grid solver: blocks versus rows (i.e strip assignment)

- Block Assignment still have a spatial locality problem on column-oriented remote data accesses
- Row-wise (strip) can perform better despite worse inherent c-to-c ratio



• Result depends on n and p

Results to show this next →

Comm-to-comp ratio: $\frac{4 \times \sqrt{p}}{n}$ for block, $\frac{2 \times p}{n}$ for strip

CMPE655 - Shaaban

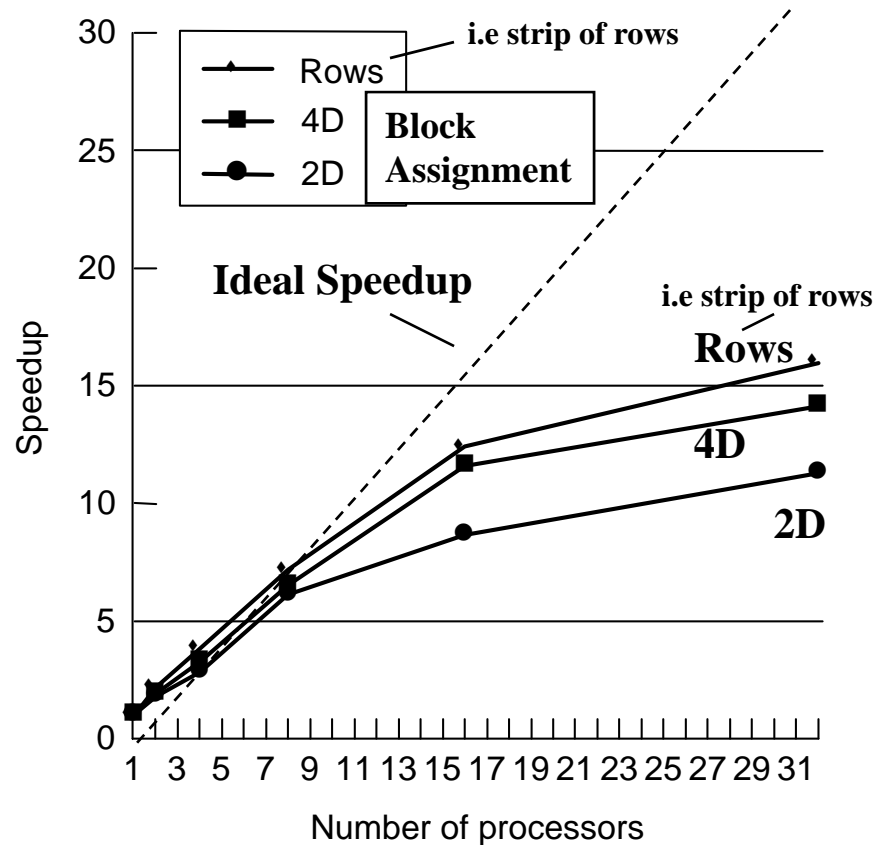
Example Performance Impact

Equation solver on SGI Origin2000 (distributed shared memory)

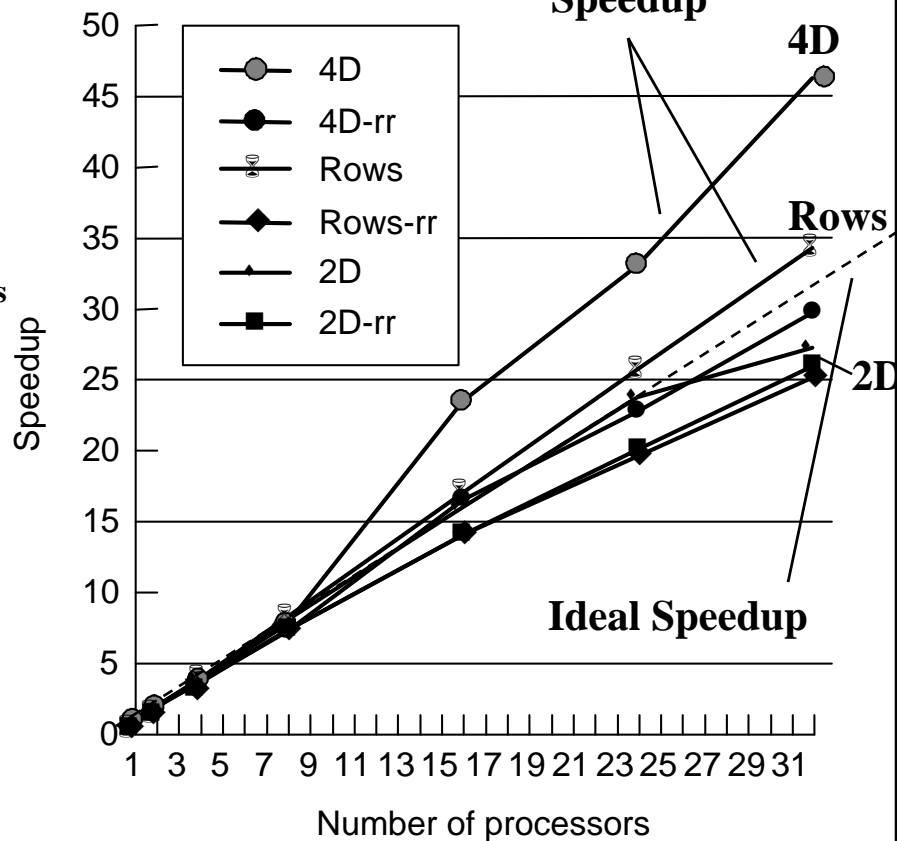
rr = Round Robin Page Distribution

Rows = Strip Assignment

Why?



514 x 514 grids



12k x 12k grids

Comm-to-comp ratio: $\frac{4 \times \sqrt{p}}{n}$ for block, $\frac{2 \times p}{n}$ for strip

CMPE655 - Shaaban

Orchestration for Performance:

Structuring Communication

Given amount of comm. (inherent or artifactual), goal is to reduce cost ↓

- Total cost of communication as seen by process:

Want to reduce

$$C = f * \left(o + l + \frac{n_c/m}{B} + t_c - \text{overlap} \right)$$

Cost of a message

Want to reduce ↓

Want to increase ↑

Latency of a message

Want to reduce ↓

Reduce ?

- f = frequency of messages
- o = overhead per message (at both ends)
- l = network delay per message
- n_c = total data sent
- m = number of messages
- B = bandwidth along path (determined by network, NI, assist)
- t_c = cost induced by contention per message
- overlap = amount of latency hidden by overlap with comp. or other comm.

n_c/m average length of message

i.e. Message Size

One may consider $m = f$

- Portion in parentheses is cost of a message (as seen by processor)
- That portion, ignoring overlap, is latency of a message
- Goal: 1- reduce terms in communication latency and
2- increase overlap

Communication Cost: Actual time added to parallel execution time as a result of communication

CMPE655 - Shaaban

Reducing Cost of Communication:

$$f * o$$

Reducing Overall Communication Overhead

i.e total

- Can reduce number of messages f or reduce overhead per message o
- Message overhead, o is usually determined by hardware and system software (implementation cost of comm. primitives)
 - Program should try to reduce number of messages m by combining messages.
i.e Fewer larger messages
 - More control when communication is explicit (message-passing).
- Combining data into larger messages: *to reduce number of messages, f*
 - Easy for regular, coarse-grained communication
 - Can be difficult for irregular, naturally fine-grained communication.
e.g duplicate computations
 - May require changes to algorithm and extra work
 - Combining data and determining what and to whom to send
 - May increase synchronization wait time.

Reduce total comm. overhead, How?

Longer synch wait to get more results data computed to send in larger message

CMPE655 - Shaaban

Reducing Cost of Communication:

Reducing Network Delay

Or Latency Total Number of Messages

- **Total network delay component** $= f * l = f * h * t_h$

Depends on
Mapping/Routing
Network Topology
Network Properties

- h = number of hops traversed in network

- t_h = link+switch latency per hop

in route from source
to destination

- **Reducing f : Communicate less, or make messages larger**
- **Reducing h (*number of hops*):**

Thus fewer messages

- Map task communication patterns to network topology
e.g. nearest-neighbor on mesh and ring etc.

- How important is this?

- Used to be a major focus of parallel algorithm design
- Depends on number of processors, how t_h , compares with other components, network topology and properties
- Less important on modern machines

Graph
Matching
Problem

Optimal
solution
is NP problem

- (e.g. Generic Parallel Machine)

Where equal communication time/delay between any two nodes is assumed (i.e symmetric network)

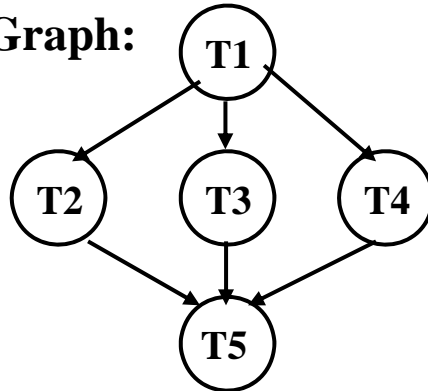
CMPE655 - Shaaban

Mapping of Task Communication Patterns to Topology

Reducing Network Delay: Reduce Number of Hops

Example

Task Graph:

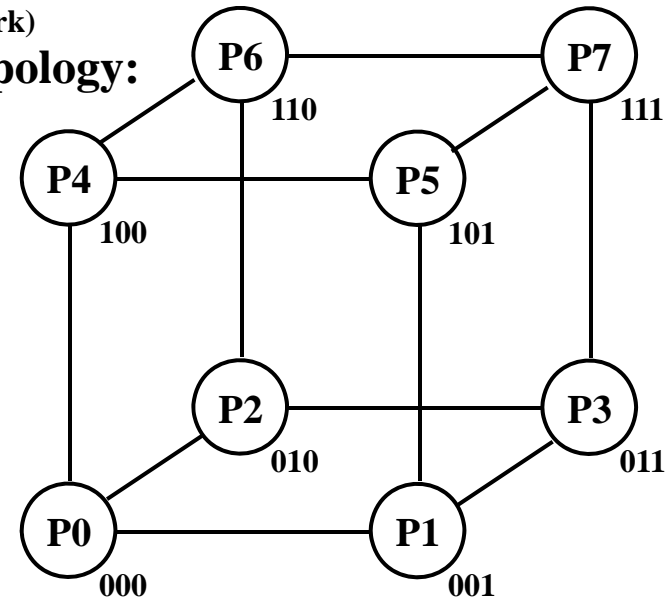


Poor Mapping:

T1 runs on P0
 T2 runs on P5
 T3 runs on P6
 T4 runs on P7
 T5 runs on P0

- Communication from T1 to T2 requires 2 hops
Route: P0-P1-P5
- Communication from T1 to T3 requires 2 hops
Route: P0-P2-P6
- Communication from T1 to T4 requires 3 hops
Route: P0-P1-P3-P7
- Communication from T2, T3, T4 to T5
 - similar routes to above reversed (2-3 hops)

(network)
 Parallel System Topology:
 3D Binary Hypercube



Better Mapping:

T1 runs on P0
 T2 runs on P1
 T3 runs on P2
 T4 runs on P4
 T5 runs on P0

- Communication between any two communicating (dependant) tasks requires just 1 hop

CMPE655 - Shaaban

Reducing Contention t_c

- All resources have nonzero occupancy (busy time):
 - Memory, communication assist (CA), network link, etc.
 - Can only handle so many transactions per unit time.
 - Contention results in queuing delays at the busy resource.
 - i.e Occupancy
 - i.e contended
- Effects of contention:
 - Increased end-to-end cost for messages. e.g delay, latency
 - Reduced available bandwidth for individual messages.
 - Causes imbalances across processors.
- Particularly insidious performance problem:
 - Easy to ignore when programming
 - Slows down messages that don't even need that resource
 - By causing other dependent resources to also congest Ripple effect
 - Effect can be devastating: *Don't flood a resource!*

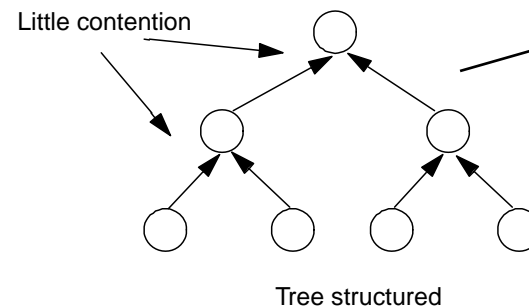
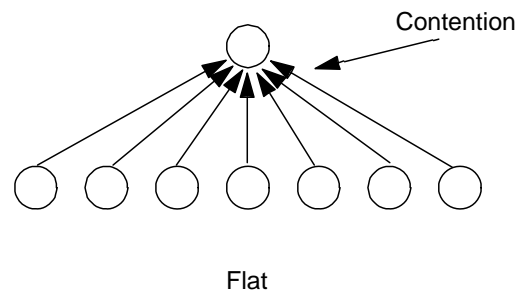
How?

CMPE655 - Shaaban

Types of Contention

- Network contention and end-point contention (*hot-spots*)
- Location and Module Hot-spots:
 - Location: e.g. accumulating into global variable, barrier
 - Possible solution: tree-structured communication

i.e one point of contention



More on this next lecture - Implementations of barriers

i.e several points of contention

- Module: all-to-all personalized comm. in matrix transpose
 - Solution: stagger access by different processors to same node temporally
- In general, reduce burstiness (smaller messages); may conflict with making messages larger (to reduce number of messages)

How to reduce contention?

Overlapping Communication

- Cannot afford to stall/wait for high latencies
- Overlap with computation or other communication to hide latency → To reduce communication cost

- Common Techniques:

Over network

- 1 – Prefetching (start access or communication before needed)
- 2 – Block data transfer (may introduce extra communication)
- 3 – Proceeding past communication (e.g. non-blocking receive)
- 4 – Multithreading (switch to another ready thread or task)

- In general these above techniques require:

i.e other ready task of the same problem

- 1 – Extra concurrency per node (*slackness*) to find some other computation.
- 2 – Higher available network bandwidth (for prefetching).
- 3 – Availability of communication primitives that support overlap.

More on these techniques in PCA Chapter 11

CMPE655 - Shaaban

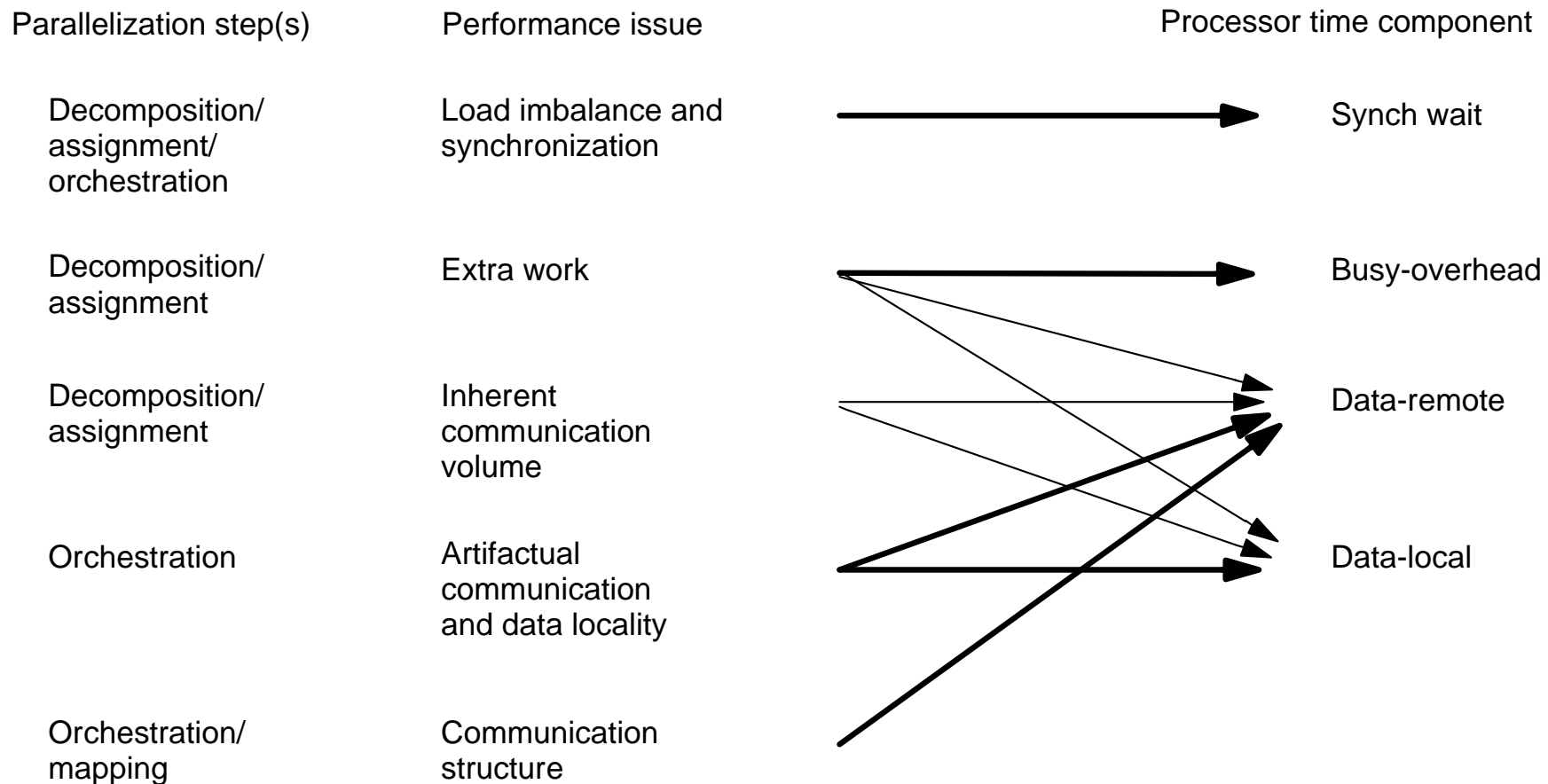
Summary of Tradeoffs

- **Different goals often have conflicting demands:**
 - **Better Load Balance Implies:**
 - Fine-grain tasks
 - Random or dynamic assignment
 - **Lower Amount of Communication Implies:**
 - Usually coarse grain tasks
 - Decompose to obtain locality: not random/dynamic
 - **Lower Extra Work Implies:**
 - Coarse grain tasks
 - Simple assignment
 - **Lower Communication Cost Implies:**
 - Big transfers: to amortize overhead and latency
 - Small transfers: to reduce contention

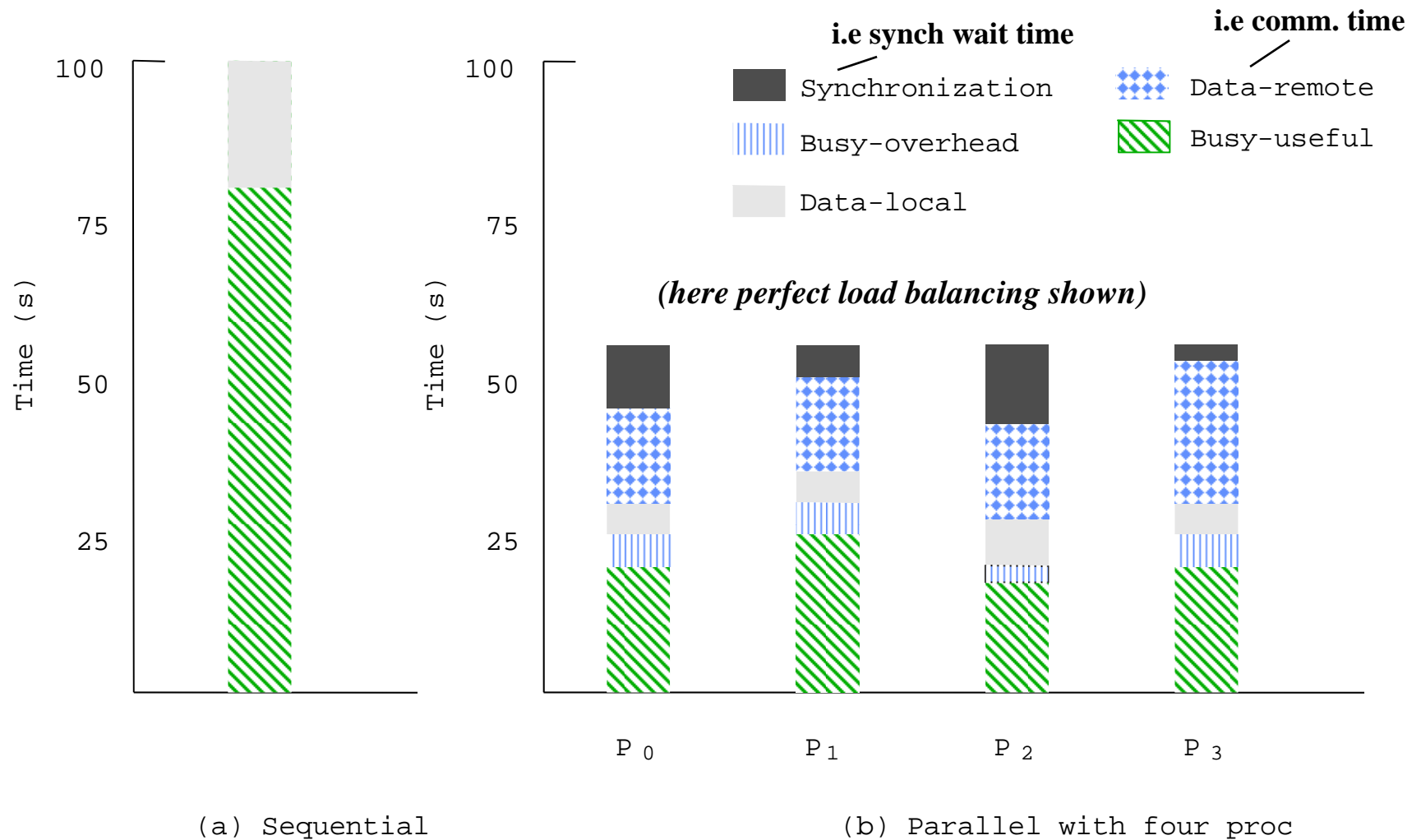
So, big or small data transfers ??

CMPE655 - Shaaban

Relationship Between Perspectives



Components of Execution Time From Processor Perspective



Summary

$$Speedup_{prob}(p) = \frac{Busy(1) + Data(1)}{\text{Max}_{(on\ any\ processor)}(Busy_{useful}(p) + Data_{local}(p) + Synch(p) + Date_{remote}(p) + Busy_{overhead}(p))}$$

- **Goal is to reduce denominator components**
- **Both programmer and system have a role to play**
- **Architecture cannot do much about load imbalance or too much communication**
- **But it can help:**
 - **Reduce incentive for creating ill-behaved programs (efficient naming, communication and synchronization)**
 - **Reduce artifactual communication**
 - **Provide efficient naming for flexible assignment**
 - **Allow effective overlapping of communication**

May introduce it , though