# Dynamic Mapping ( dynamic load balancing)

Centralized ← → Distributed

## Centralized

→ All executable Tasks ⎯⎯ Common Central Data Structure
                           ⎯⎯ special process (Master + slaves)
                           ⎯⎯ subset of processes

→ A process with no work takes portion of available work from the central data structure / Master

→ New tasks are added to Central data structure / reported to the Master

→ Easier to implement than distributed

→ Limited Scalability

e.g. Quick Sort each row of an $n \times n$ matrix

            ↙  ↘

  $O(n)$        $O(n \log n)$
   ( ? )

⇒ Equal number of rows to each process - imbalance

⇒ Maintain Central pool of indices of rows yet to be sorted
   Process when idle can pick an index from this pool   ( self - scheduling)

→ Good load balancing but bottlenecks in accessing shared work queue if each task does not have large enough computation

∴ To ease the bottleneck, assign more than one Task at a time ⟹ Chunk scheduling

But what if there is a load-imbalance due to large chunk sizes?

⟹ Reduce chunk size as the program progresses
(linear, non-linear, ...)

## Distributed

→ Set of executable Tasks are distributed among Processes

→ Processes exchange tasks at run time to balance work

→ No bottleneck of centralized schemes

Its critical parameters are:

→ Pairing of Sender/receiver processes
→ Initiation of work transfer - Sender or Receiver?
→ Amount of work transfer
→ Point of work transfer

To be effective in Message passing platforms, Size of Tasks (Computations) should be larger than size of data associated with it.

# Reduce Process Interaction Overheads

→ Volume of data

→ Frequency of Interaction

→ Spatial & Temporal pattern of Interactions

Techniques to Reduce Interaction Overheads:

→ Promote use of recently fetched data (Temporal locality)

→ Minimize the Overall volume of shared data that needs to be accessed by concurrent processes

$$\frac{2n^2}{p} \quad Vs \quad \frac{n^2}{\sqrt{p}} + n^2$$

$$(1-d) \qquad\qquad (2-d)$$

→ Locally Compute, Store intermediate results, perform shared data access only to place final results

```
for(i = 0; i < N; i++)
    dp += A[i] * B[i];
```
---
```
for( i = id * N/p ; i < (id+1)N/p ; i++)
    ldp +=  A[i] * B[i];

dp += ldp;
```

→ Minimize frequency of interactions as each interaction involves Relatively high Startup cost

- Restructure the algorithm such that shared data are accessed and used in large pieces

Amortizes startup cost over multiple accesses (akin to large cache lines for spatial locality)

Leads to fewer messages on Message Passing systems

e.g. Sparse Matrix Vector Multiplication

P0 needs $b[3]$, $b[5]$ & $b[7]$, say

$\downarrow P_1 \downarrow$  $\downarrow$ P2

Combine interaction with P1 to get $b[3]$ & $b[5]$ instead of 2 separate interactions.

→ Minimize contention

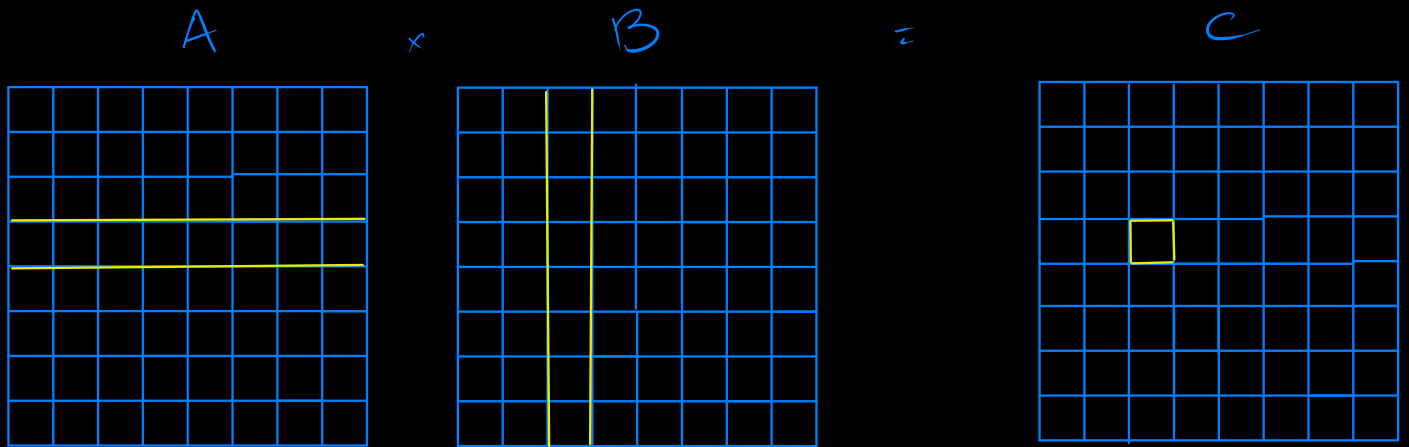Why? Multiple Tasks accessing same resource Concurrently

Only one of the operations can proceed at a time, others are queued, proceed sequentially

e.g. Multiple simultaneous transmission of data over the same interconnect link

Multiple simultaneous access to same memory block

Multiple processes sending messages to the same process at the same time

$C = AB$ using 2-d Partitioning, p tasks, p processes

(1-1 mapping)

∴ Each Task to Compute $C_{i,j}$, $0 \le i, j < \sqrt{p}$

Let $p = 64$, $c_{ij} = \sum_{k=0}^{\sqrt{p}-1} A_{ik} * B_{kj}$

A     x     B     =     C



Processes Computing $C_{0,0}, C_{0,1}, \dots, C_{0,\sqrt{p}-1}$ will attempt to read $A_{0,0}$ at the same time.

Processes Computing $C_{0,0}, C_{1,0}, \dots, C_{\sqrt{p}-1,0}$ will attempt to read $B_{0,0}$ at the same time.

⇒ Contention in NUMA shared address space

      & Message passing platforms

∴ Redesign the algorithm to access data in contention-free patterns.

$$C_{i,j} = \sum_{k=0}^{\sqrt{p}-1} A_{i, (i+j+k)\%\sqrt{p}} * B_{(i+j+k)\%\sqrt{p}, j}$$

Process Computing $c_{0,0}$ accesses: $A_{0,0}$

    ʼʼ     ʼʼ    $C_{0,1}$    ʼʼ    :   $A_{0,1}$

           ⋮

         $C_{0,7}$     ʼʼ    :   $A_{0,7}$

Process computing $C_{5,0}$ accesses: $A_{5,5}$

$\quad\quad\quad\quad\quad C_{5,1}\quad\quad : \quad A_{5,6}$

$\quad\quad\quad\quad\quad C_{5,2}\quad\quad : \quad A_{5,7}$

$\quad\quad\quad\quad\quad C_{5,3}\quad\quad : \quad A_{5,0}$

$\quad\quad\quad\quad\quad C_{5,4}\quad\quad : \quad A_{5,1}$

$\quad\quad\quad\quad\quad C_{5,5}\quad\quad : \quad A_{5,2}$

$\quad\quad\quad\quad\quad C_{5,6}\quad\quad : \quad A_{5,3}$

$\quad\quad\quad\quad\quad C_{5,7}\quad\quad : \quad A_{5,4}$

[ Centralized schemes in Dynamic Mapping are frequent source of contention for shared data structures ⇒ choose distributed mapping ]

→ Overlap Computations with Interactions

- Perform useful computations while waiting for shared data to arrive after initiating the interaction.

   - Initiate interaction early enough by identifying computations that can be performed before the interaction.

   - If one Task blocks on an interaction, and if there is another Task available to the same process, execute the other Task.

- Initiate advance work transfer if the process can anticipate running out of current work.

    - Requires support from the Programming Paradigm and underlying hardware

    e.g. Non-Blocking Message Passing Primitives
    Prefetching hardware/compiler to advance loads

→ Replicating Data or Computations

    - For read-only shared data access, replicate it in each process, thereby avoiding interaction
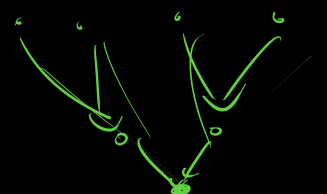
    - but it increases memory requirement linearly with number of concurrent processes.

    - Instead of sharing intermediate results, each process can compute them.

    - Trade-off between interaction overhead to replicated computation

→ Optimized Collective Operations
    - Static and regular pattern of interaction between a group of tasks e.g Broadcast
                 Reduce

- Used for → sharing data
  → communication intensive computations
  → Synchronization

- Optimized implementation of collectives are available in libraries from parallel platform vendors

- In exceptional scenarios, develop one's own collective implementation.

→ Overlap Interactions with Other Interactions

- with hardware support, overlapping interactions between multiple pairs of processes can reduce effective volume of communication.

e.g. One to all broadcast from $P_0$ to $P_1, P_2 \& P_3$

Naive algorithm:

$$P_0 \xrightarrow{1} P_1 \xrightarrow{2} P_2 \xrightarrow{3} P_3 = 3 \text{ units of time}$$

Overlapped:

$$P_0 \longrightarrow P_1 \qquad P_0 \longrightarrow P_2$$
$$P_1 \longrightarrow P_3$$
$$\qquad\qquad 1 \qquad\qquad\quad 2 \qquad\qquad = 2 \text{ units of time}$$

what if 4 messages need to be sent?

$$P_0 \xrightarrow{m_0} P_1 \quad P_0 \xrightarrow{m_0} P_2 \qquad\qquad P_0 \xrightarrow{m_1} P_1 \quad P_0 \xrightarrow{m_1} P_2$$
$$P_1 \xrightarrow{m_0} P_3 \qquad\qquad\qquad\qquad P_1 \xrightarrow{m_1} P_3$$

$$P_0 \xrightarrow{m_2} P_1, \quad P_0 \xrightarrow{m_2} P_2 \atop P_1 \xrightarrow{m_2} P_3 \qquad\qquad P_0 \xrightarrow{m_3} P_1, \quad P_0 \xrightarrow{m_3} P_2 \atop P_1 \xrightarrow{m_3} P_3 \quad = 8 \text{ units}$$
$$\text{of time}$$

If the naive algorithm is used:

$$P_0 \xrightarrow{m_0} P_1, \quad P_0 \xrightarrow{m_1} P_1, \quad P_0 \xrightarrow{m_2} P_1, \quad P_0 \xrightarrow{m_3} P_1, \quad P_1 \xrightarrow{m_3} P_2 \quad P_2 \xrightarrow{m_3} P_3$$

$$P_1 \xrightarrow{m_0} P_2 \quad P_1 \xrightarrow{m_1} P_2 \quad P_1 \xrightarrow{m_2} P_2 \quad P_2 \xrightarrow{m_2} P_3$$

$$P_2 \xrightarrow{m_0} P_3 \quad P_2 \xrightarrow{m_1} P_3$$

$$= 6 \text{ units}$$
$$\text{of time !}$$

- An exceptional scenario for a user developed customized collective operation.