

① Parallelization Abstraction

1. Parallelization steps:

1. Decomposition of applications into functional tasks/data blocks

• Factors:

1. Application Type
2. Concurrency degree (Embarrassingly parallel, Pleasingly parallel, Embar. serial)
3. Granularity (Fine, medium, coarse)
4. Target system (Shared-memory, Distributed-memory architectures)

• Types

1. Data decomposition (work sharing)
matrix mult.
2. Functional decomposition (work stealing)
diff. execution, list scheduling, work sharing, work stealing
3. Recursive decomp (Fibonacci, Merge sort)
rec. tree, pattern matching
4. Exploratory decomp
5. Speculatory (if)

• Can be combined.

2. Dependence analysis of decomposed functional task or data blocks.

Types:

1. Control dependencies (

if $u=0$ then
...
else
...

2. Data dependencies

$\begin{matrix} a:=2 \\ u:=a+2 \\ v:=a+3 \\ x:=u+v \end{matrix}$

Flow dependencies - one task writes and other reads.

3. Name dependencies

• Anti-dep - one reads and another writes a var

• Output-dep - both tasks write that var.

3. Mapping of the execution of functional task on data blocks onto target system.

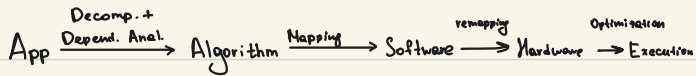
When? • Temporal ordering - ^{when} start time to each subtask.

Where? • Spatial assignment - ^{where} placement on processing elements, according to temporal ordering

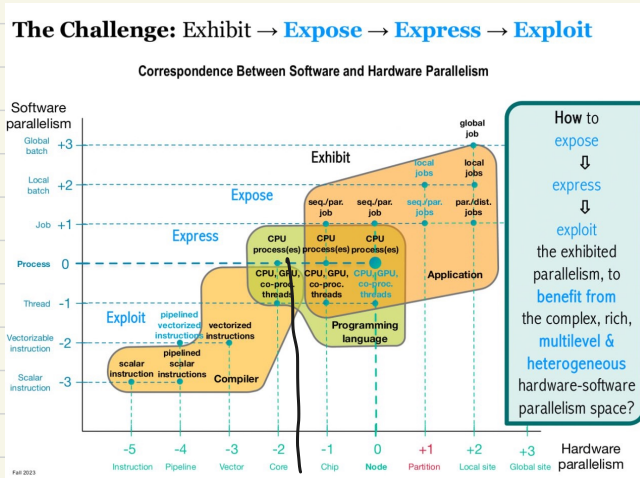
How do map? • Static | Dynamic \leftrightarrow offline | online.

4. Programming: express all the parallelism 1-3 via a prog language.

1. Shared memory (with threads)
2. Shared memory (without threads)
3. Distributed-memory
4. Data parallel
4. Hybrid distributed Shared-memory
5. Single Program Multiple Data
6. Multiple Program Multiple Data

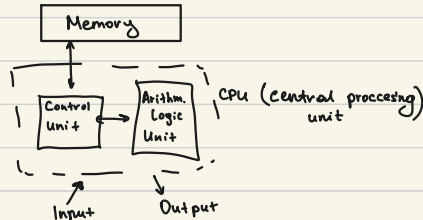


2. Parallelism: Forms and Levels

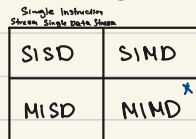


3. Architecture computer:

1. Von Neumann 1945



2. Flynn's Taxonomy (Parallel computers)



- Small processing unit (PU)
- CU, control unit

* most popular

Global task executed in parallel require synchronization/data exchange over network.

a) SIMD



b) MIMD:

Limitations: nothing about organization of memory

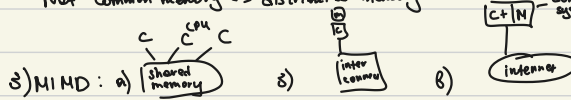
3. Flynn - Johnson's 1988

Location: Centralized / Distributed

Access policy: Full / No access

Common memory \Rightarrow shared-memory machines

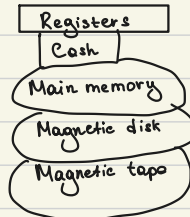
Not common memory \Rightarrow distributed-memory or message-passing machines



4.

Memory architectures

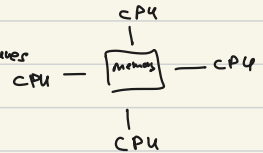
1. Hierarchy



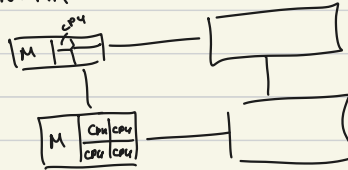
• Shared memory machines

• UMA, Uniform Memory Access architecture

• visible to all CPUs

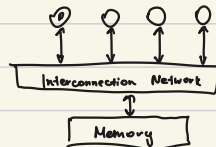
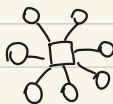


• NUMA



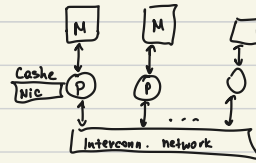
Model 1: Shared memory: (Symmetric multiprocessing and UMA)

\Rightarrow VAS + UMA



UMA + VAS (Unified address space)

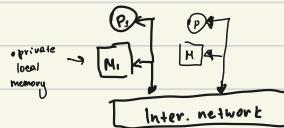
• Model 2: Distributed memory (Asymmetric multiprocessing and NUMA)



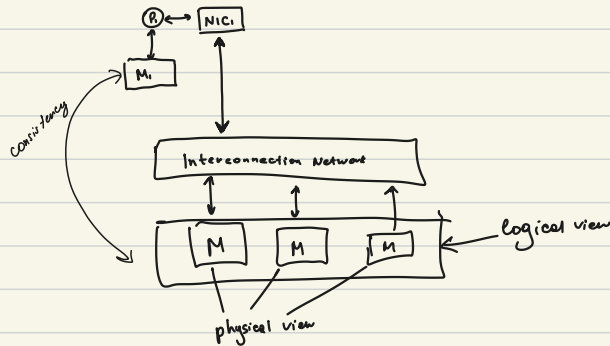
• Model 3: Hybrid distributed shared-memory (

• GAS - Global Address Space

• A/symmetric Multiprocessing



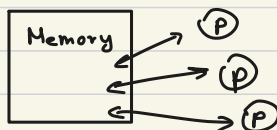
• Model 4: Distributed shared-memory



Model 5: Fully distributed (connection via WAN)

Parallel Programming models

1. Shared memory model (without threads)



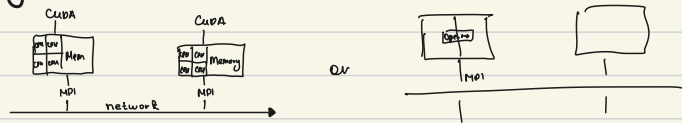
(read and write
asynchronously)

2. Shared-memory (with threads)

- Processes have threads, that can be concurrent.
- Each thread has local data, shares global memory.
- „Communicate“ via global memory.

3. Distributed-memory / Message-Passing

4. Hybrid Model



a) with MPI and CUDA

b) with MPI and OpenMP

5. SPMD

6. MPMD

Network topologies

- 1) On chip network
- 2) System / Storage area networks
- 3) Local area networks
- 4) Wide area network.

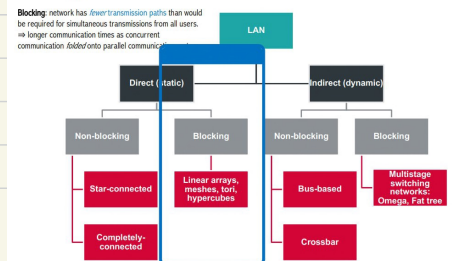
$$M = O \left(\frac{E u(d)}{2 \Delta u} \right) \Delta u =$$

$$\sigma^2 = \frac{2s^2 \ln(1.35)}{\sigma}$$

$$\epsilon^2$$

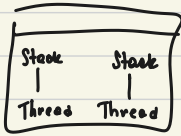
3: LAN

Blocking network has fewer transmission paths than would be required for simultaneous transmissions from all users
 \Rightarrow longer communication times as concurrent communication folds onto parallel communication



Parallel Programming Shared Memory

OpenMP



Shared-Memory
Ap

MPI



Single-thread
process

Message passing approach

